LEARNING THE HAMILTONIAN OF DISORDERED MATERIALS WITH EQUIVARIANT GRAPH NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph neural networks (GNNs) have shown promise in learning the ground-state electronic properties of molecules and crystalline materials, subverting computationally intensive density functional theory (DFT) calculations. Materials with structural disorder, however, are more challenging to learn as they exhibit higher complexity and a more extensive palette of local atomic environments, all of which require large (10+ Å) cells to be accurately captured. In this work, we adapt efficient equivariant GNN approaches to learn disordered materials' electronic properties, represented by the Hamiltonian matrix (H). Since creating a large graph corresponding to the whole structure of interest would be computationally prohibitive, we introduce an *augmented partitioning* approach in which the graph is sliced into multiple partitions, each augmented with masked virtual nodes and edges. This method maintains correct atomic neighborhoods within a single message passing layer, allowing for the network to learn the electronic properties of amorphous HfO₂ materials with 3,000 nodes (atoms), 500,000+ edges, and ~28 million orbital interactions (non-zero entries of H).

004

010 011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

Predicting structure-property relationships in atomically resolved materials lends itself optimally
to graph-structured data. Graph Neural Networks (GNNs) have proven capable of learning these
relationships while constrained by their underlying symmetries (Veličković et al. (2018)). Trained
networks have been able to learn molecular- and atomic-level quantities (Wang et al. (2023)), circumventing otherwise computationally prohibitive simulations at the *ab initio* level. More recently,
GNNs have been adapted to predict electronic properties, described by a discretized ground-state
Hamiltonian matrix *H*. Rapid and accurate constructions of *H* can unlock *in silico* explorations of
the large design space of electronic materials (Klinkert et al. (2020)).

The matrix H can be decomposed into sub-matrices $H_{i,j}$ that encode the coupling between the sets of atomic orbitals located on atoms i and j. These coupling terms are a function of the identity and relative coordinates of the local environment. Predicting such electronically resolved information introduces additional challenges over atomically resolved quantities, as the output data is equivariant under rotation. Existing work on electronic property prediction has mainly treated the cases of small molecules (Zhong et al. (2023); Yu et al. (2023b); Bai et al. (2021)) and well-ordered materials (Li et al. (2022); Gong et al. (2023); Wang et al. (2024a)). The graph representations of these structures are fairly small - small molecules contain only a few atoms, and in crystalline materials all relevant structural information can be captured within the smallest repeating unit cell.

Many applications, however, require the computation of the electronic properties of materials with structural disorder, such as local or extended defects, or in an amorphous phase (Ducry et al. (2020); Kaniselvan et al. (2023); Strand et al. (2018)). These materials typically contain a limited number of atomic species, from 2 to 5, with quasi-random distributions. Accurate *ab initio* simulations of the resulting disordered atomic structures are only possible if large unit cells composed of hundreds to thousands of atoms are used (Repa & Fredin (2023)). As a consequence, prohibitively expensive computations must be performed with a density functional theory (DFT) tool to obtain their electronic properties. The prospect of applying deep-learning solutions to handle such materials is thus particularly attractive. To be of practical relevance, however, they should be able to generalize to large scales.

Here, we extend equivariant GNN approaches to learn and predict the electronic properties of materials in amorphous phases by fitting the sub-blocks of the ground-state Hamiltonian *H* in matrix form. Our main contributions are:

- We develop an efficient GNN-based model for electronic property prediction, by combining (1) the SO(2)-convolution approach detailed in Passaro & Zitnick (2023), (2) the equivariant attention mechanism introduced in Liao et al. (2023), and (3) concepts from Gong et al. (2023) and Wang et al. (2024a) to introduce learnable node/edge embeddings along with basis transformation layer to pre-process the targets and map predictions to the Hamiltonian output. We provide the code for this implementation in [the Supporting Materials].
- We propose an efficient augmented partitioning method that breaks down input graphs into small pieces and corrects atomic environments with masked virtual nodes and edges. This allows arbitrarily large graphs to be decomposed into independent partitions that can fit into GPU memory during training without compromising the achievable testing accuracy. Our approach enables the training and prediction of unfeasibly large systems including realistic amorphous materials and heterostructures that can contain up to hundreds of thousands of atoms in a unit cell.

071 We combine our model and *augmented partitioning* approach to treat a real example with practical scientific relevance. Specifically, we consider hafnium dioxide (HfO_2), one of the most technologically relevant amorphous oxides (Choi et al. (2011)). The theoretical study of defects and transport 073 properties in HfO2 is relevant to several research areas, from optimizing gate dielectrics for transis-074 tors (Strand et al. (2018)) to developing new resistive-switching technologies enabling in-memory 075 computing (Kaniselvan et al. (2023)). With this, we achieve a prediction accuracy of 5.87 meV, 076 matching the eigenvalues of H to within 0.87% relative L1 error, on structures with 3,000 atoms, 077 which require several (3.65) hours to compute using DFT. Our work advances applications of equiv-078 ariant GNNs towards practical use cases in computational physics, chemistry, and materials science. 079

080 081

058

060

061

062

063

064

065

067

068

069

2 BACKGROUND & RELATED WORK

082 The electronic properties of a material refer to its set of 083 energy levels (ε) and wavefunctions (ψ) that electrons 084 can occupy. They correspond to the eigenvalues and 085 eigenvectors of the Hamiltonian matrix H describing the atomic system of interest. This quantity is a function of 087 the location (relative positions $\{r_i\}$) and identity (atomic numbers $\{Z_i\}$) of all constituent atoms $\{i\}$ (Hohenberg & Kohn (1964)). Therefore, predicting the electronic properties consists of learning the mapping $F : \{r_i, Z_i\} \rightarrow$ 090 H between the atomic structure and the elements of the 091 corresponding Hamiltonian matrix (Fig. 1). 092

The entries of the ground-state Hamiltonian matrix H094 are typically computed from first-principles with DFT 095 (Kohn & Sham (1965)). In several widely used codes, the wavefunctions are expanded into a basis $|\varphi\rangle$ of non-096 orthogonal atomic orbitals localized around atomic positions, each built, for example, from contracted Gaussian 098 functions (Kühne et al. (2020); Neese (2011)). These orbitals transform like spherical harmonics under rotation 100 $\hat{\boldsymbol{r}}
ightarrow \hat{\boldsymbol{r}}'$: $Y_m^l(\hat{\boldsymbol{r}}') = \hat{\sum}_{m'} \boldsymbol{D}_{mm'}^l(\boldsymbol{R}) Y_{m'}^l(\hat{\boldsymbol{r}})$. Here, Y_m^l 101 is the spherical harmonic of degree l and order $m \in$ 102



Figure 1: Schematic of the mapping between the atomic graph and the blocks of the Hamiltonian matrix H in the localized orbital basis of choice. Each orbital block represents the couplings between atomic orbitals on the same atom ($H_{i,i}$, diagonal) or between two different atoms within r_{cut} ($H_{i,j}$, off-diagonal).

 $\{-l, \dots, l\}. D_{mm'}^{l}(\mathbf{R}) \text{ is the Wigner-D matrix of degree } l \text{ corresponding to the rotation } \mathbf{R}, \text{ which} \\ \text{transforms the corresponding spherical harmonic. } \hat{r} \text{ and } \hat{r}' \text{ are normalized direction vectors.} \end{cases}$

The localized nature of the basis states leads to finite spatial overlaps between them. The resulting Schrödinger equation at the core of DFT takes the form of a generalized eigenvalue problem: $H\psi = \varepsilon S\psi$. Here, the Hamiltonian matrix $H^{(N\times N)}$ has entries $H_{i,j} = \langle \varphi_i | \hat{H}(\mathbf{r}) | \varphi_j \rangle$ where $\hat{H}(\mathbf{r})$ is the so-called Hamiltonian operator, while the Overlap matrix $S^{(N\times N)}$ is made of $S_{i,j} = \langle \varphi_i | \varphi_j \rangle$. They are both coarse-grained matrices of size $N = \sum_{k} N_{atoms}^{k} \cdot N_{orb}^{k}$, where N_{atoms} is the number of atoms, N_{orb} the number of orbitals per atom, and k indexes over the different atomic species found in the system. Note that S reduces to the identity matrix in case of an orthogonal basis $|\varphi\rangle$. Otherwise, it can be directly computed from the basis as the problem's physics does not influence it.

The Hamiltonian matrix can be decomposed into sub-matrices $H_{i,j}$ of size $(N_{orb}^i \times N_{orb}^j)$, each describing the interactions between all basis elements (orbitals) on atoms *i* and *j*. Diagonal blocks ($H_{i,i}$) are the interactions between orbitals on the same atom. When represented on a local basis, the matrix is near-sighted; the interactions between orbitals on different atoms decay exponentially with increasing interatomic distance. Since an atomic orbital basis is used, the sub-matrices are equivariant under rotation of the atomic bonds, with their transformation properties related by the Wigner-D matrix.

119 120

121

2.1 CHALLENGES UNIQUE TO DISORDERED MATERIALS

122 Computing the electronic properties of disordered materials with DFT still re-123 quires defining a repeating 'unit cell' and 124 using periodic boundaries to avoid dan-125 gling bonds. This periodicity, however, 126 can alter the material's amorphous na-127 ture if the repeating unit cell is too small. 128 Atoms can interact with all their peri-129 odic images, leading to non-physical phe-130 nomena such as the formation of coherent 131 electronic states across cells. These phe-132 nomena can be prevented by constructing 133 'large-enough' unit cells (12+ Å (Repa & Fredin (2023)) to a few nanometers 134 (Ducry et al. (2020))) that better approx-135 imate disorder. Generating the Hamilto-136



Figure 2: Illustration of the difference between ordered (left) and disordered materials (right), whose structural features can only be captured by defining a large unit cell. In both cases, the smallest repeating unit cell is delimited by a black box, while the circles/lines correspond to atoms/bonds.

nian matrix H of these systems with DFT involves tens to hundreds of self-consistent field (SCF) loops, each requiring the diagonalization of an intermediate H. As this numerical operation scales with $\mathcal{O}(N_{atoms}^3)$, analyzing electronic properties for large amorphous systems (or different amorphous representations of the same material) is very often computationally unaffordable.

- 140 141
- 142

2.2 DEVELOPMENT OF MODELS FOR THE PREDICTION OF ELECTRONIC PROPERTIES

¹⁴³ Only a few studies have attempted to directly predict the Hamiltonian matrix H rather than directly ¹⁴⁴ fitting invariant quantities such as the total energy. The key is to constrain the solution space by ¹⁴⁵ leveraging prior knowledge of physical symmetries, e.g., rotational equivariance of orbital blocks.

146 Several works leverage local structure descriptors containing sufficient detail about each atom's en-147 vironment. The mapping between these descriptors and the orbital blocks of the Hamiltonian was 148 then learned through Kernel Ridge Regression (Hegde & Bowen (2017)) or multilayer neural net-149 works (Schütt et al. (2019); Gu et al. (2024)). Certain descriptors, such as Atomic Cluster Expansion 150 (ACE) (Drautz (2020)), can be extended to orbitally resolved data (Nigam et al. (2022)) and used 151 to predict the Hamiltonian matrix of crystalline materials (Zhang et al. (2022)). These approaches 152 achieved prediction accuracy of a few meV on small-molecule datasets. Rotational equivariance was 153 enforced through the descriptors (Zhang et al. (2022)) or data augmentation (Schütt et al. (2019)). Initial GNN-based approaches, such as the network developed by Li et al. (2022), are intrinsically 154 invariant to the translation and permutation of the inputs. Information about the rotational equivari-155 ance was incorporated by rotating to a pre-selected axis before training, which reduces the problem 156 to a rotationally invariant one. 157

In equivariant GNNs, the predicted Hamiltonian rotates along with the input (Yu et al. (2023b); Zhang et al. (2024); Batatia et al. (2023); Gong et al. (2023)), which requires maintaining SO(3)equivariance within the model. This means that all network operations f acting on input embedding x^{l} of degree l must satisfy: $f(D^{l}(R) \cdot x^{l}) = D^{l}(R) \cdot f(x^{l})$. The networks are trained using Message Passing (MP), where each MP layer works as follows: An atom i receives input messages

162 from each neighboring source atom j. Each input message goes through convolution operations 163 that combine features with different l while preserving equivariance; a specific output embedding 164 $x_{ji}^{l_3}$ of degree l_3 can be computed through: $x_{ji}^{l_3} = \sum_{l_1, l_2} x_j^{l_1} \otimes_{l_1, l_2}^{l_3} h_{l_1, l_2, l_3} Y^{l_2}(\hat{r}_{ji})$. Here, \hat{r}_{ji} is a normalized vector indicating the direction of the edge connecting the atoms j and i, and h is a set 165 166 of trainable weights. The sum runs over tensor products which take x_i (a source input embedding 167 of degree l_1) and Y^{l_2} (a filter spherical harmonic embedding of degree l_2) and produce the output 168 embedding:

169
170
$$(\boldsymbol{x}_{j}^{l_{1}} \otimes_{l_{1},l_{2}}^{l_{3}} Y^{l_{2}}(\hat{\boldsymbol{r}}_{ji}))_{m_{3}}^{l_{3}} = \sum_{m_{1},m_{2}} (\boldsymbol{x}_{j}^{l_{1}})_{m_{1}} C_{(l_{1},m_{1}),(l_{2},m_{2})}^{l_{3},m_{3}} h_{l_{1},l_{2},l_{3}} Y_{m_{2}}^{l_{2}}(\hat{\boldsymbol{r}}_{ji}),$$

171

where the $C^{l_3,m_3}_{(l_1,m_1),(l_2,m_2)}$ are the Clebsch-Gordan coefficients that are indexed by the order m and 172 degree l of the input, filter, and output embeddings. The combination of feature (x) and geometric 173 (\hat{r}) information along each edge encodes both the identity and structure of the system. These 'Tensor 174 Field Networks' (TFNs) (Thomas et al. (2018)) achieve state-of-the-art accuracy on small molecule 175 (Yu et al. (2023b)) and crystalline (Gong et al. (2023)) datasets. However, they are also much more 176 computationally expensive. The network training scales with $O(l_{max}^6)$, where l_{max} is the maximum 177 degree of the angular momentum considered. Fully, E(3)-equivariant networks are difficult to apply 178 beyond a few atoms (Zhang et al. (2024)).

179 Recently, the computational cost of training equivariant GNNs has been significantly reduced by 180 combining the benefits of data rotation and equivariant network operations. These approaches take 181 advantage of the fact that when edges are rotated to align with a fixed axis (y or z depending on con-182 vention), the only non-zero spherical harmonic components are those of order m = 0. By keeping 183 track of the bond vectors and performing internal spherical rotations, complex SO(3) convolutions 184 can thus be reduced into SO(2) linear convolution operations (Passaro & Zitnick (2023)). Addition-185 ally, under these conditions, the Clebsch-Gordan coefficients exhibit a predictable sparsity pattern (non-zero only when $m_3 = \pm m_1$). Altogether, the scaling reduces to $O(l_{max}^3)$ (Wang et al. (2024a)), 186 and the network's training speeds up, enabling the use of higher-order angular momenta (l_{max}) and 187 more parameters to capture finer, more complex details of the surrounding environment. An ad-188 vanced SO(2) convolution network was developed by Passaro & Zitnick (2023) (eSCN)) and further 189 expanded by Liao et al. (2023) (EquiformerV2) with the inclusion of equivariant attention and sep-190 arable activation layers. A subsequent implementation of this approach on Hamiltonians by Wang 191 et al. (2024a) achieved better performance on custom crystalline 2D-material datasets compared to 192 previous tensor field and invariant networks. 193

3 **METHODS**

We adapt the 'EquiformerV2' network by incorporating concepts from Gong et al. (2023) and Wang et al. (2024a). In this section, we present an overview of the methods used to initialize the graph, construct the network, and propose an efficient augmented partitioning approach to train it. Relevant implementation details and ablation studies are presented in Section 4 and Appendices A, C, D.

3.1 NETWORK LAYOUT

202

194

195 196

197

199

200 201

203



209 Figure 3: High-level overview of the network, illustrating the update of node features (n_X) and edge features 210 (e_X) after message passing layer X (MP X). Z refers to atomic numbers, while r_{ij} is the set of scalar distances between atoms. 211

212

213 When constructing a given material's graph, we can leverage the Hamiltonian's near-sightedness to only retain edges corresponding to atoms within a fixed interaction distance r_{cut} , beyond which 214 the interactions are negligible. Orbitals located on different atoms can interact with each other over 215 distances of ~ 10 Å, giving rise to specific off-diagonal blocks in the Hamiltonian matrix.

216 The graph's nodes/edges are initialized with an embedding of shape $(N_n, (l_{max} + 1)^2)$ 217 $d_{sphere,n}$ //(N_e , $(l_{max} + 1)^2$, $d_{sphere,e}$), where $N_{n/e}$ is the number of nodes/edges, $d_{sphere,n/e}$ 218 the channel dimension for embeddings, and l_{max} the maximum degree of the features. The l = 0219 channels of the node embeddings are initialized with atomic numbers, while the l = 0 of the edge 220 embeddings are initialized with the scalar distance between the two connecting nodes, expanded in the chosen basis, here Gaussian functions. All other components are initially set to 0. The orbital 221 blocks representing the interaction between atoms are flattened into 1D tensors to form the labels 222 for each node/edge during the supervised training process. They are then converted from uncoupled 223 to coupled basis using a Wigner-Eckart transformation (Appendix A.1). 224

225 During training, each MP layer updates the node, and then the edge representations. Within the 226 node update block, each node *i* receives messages from each of its neighbors j, consisting of the concatenated embeddings n_i , n_j , and that of the edge e_{ii} , rotated to align with the z-axis. After 227 SO(2) convolutions are performed on the input messages, the resulting output messages are rotated 228 back to their original orientations and aggregated onto the node i to update its embedding n_i . The 229 updated node embeddings are then used to update the edges through a similar process, without the 230 attention layer. Subsequently, the node embeddings in the next message-passing layer are updated 231 with the new edge embeddings. A high-level overview of the network is presented in Fig. 3, and 232 a more detailed outline can be found in Appendix A. During the inference phase, the output is 233 converted back to the uncoupled basis through a Wigner-Eckart layer (Appendix A.1) to reconstruct 234 the Hamiltonian matrix **H**. 235

Our network architecture is similar to that of EquiformerV2 and DeepH2. A particular difference lies in the use of gate activation layers instead of the S² activation layer of EquiformerV2. Although the S² layer allows for non-linearity to be introduced to higher-order features, achieving numerically perfect equivariance requires a large grid size, making it computationally expensive when implemented on large graphs (Passaro & Zitnick (2023)).

3.2 AUGMENTED PARTITIONING WITH MASKED VIRTUAL NODES/EDGES

Training GNN representations of the large amorphous materials considered here incurs
high memory consumption and long computation times per epoch (Appendix G.2). The dense connectivity of these graphs also leads to heavy communication overhead in full-batch distributed approaches (Wan et al. (2022)).

241

242

250 A variety of methods have been developed to reduce the memory requirements and increase 251 the amount of parallelism during training of 252 large graphs (Besta & Hoefler (2024)). How-253 ever, they do not cover the specificity of our 254 application. Much of previous work is con-255 cerned with modifying connectivity to effec-256 tively propagate long-range information, for ex-257 ample, by dividing the graph into sub-graphs 258 and passing messages intra- and inter-subgraph 259 (Liao et al. (2018)), or introducing virtual 260 nodes to transfer messages over longer dis-261 tances (Qian et al. (2024)). Such approaches are unnecessary in our case since the graph is near-262 sighted and long-range interactions are negligi-263 ble (Appendix F.1). 264



Figure 4: Illustration of how connectivity beyond the partition boundaries is incorporated, using oneway masked virtual edges $(e_{j' \rightarrow i})$ from masked virtual nodes (j') to the labeled nodes. Edges between nodes within the partition are two-way $(e_{j\leftrightarrow i})$. A set of edges from labeled and virtual nodes is shown for a single node (i) within the partition. The solid vertical lines are the partition boundaries. Different colors represent different atomic species.

Neighborhood sampling techniques such as ClusterGCN (Chiang et al. (2019)) represent an alter native, but they sacrifice the exact connectivity of the graph to enable computational efficiency and
 scalability. As *H* is a function of the relative positions and species of all constituent atoms, directly
 changing the representation of local atomic neighborhoods compromises the achievable accuracy.
 Omitting any connections leads to misinformation and poor generalization, as the network tries to
 fit to the target data while aggregating information through an incorrect/incomplete graph structure.

Training of large GNNs for electronic property prediction thus necessitates new strategies to become computationally feasible.

To enable efficient training of the graph while maintaining correct atomic environments and neigh-273 boring edge connections, we introduce an *augmented partitioning* approach. A visual representation 274 of it is provided in **Fig. 4**. The graph is partitioned into slices along the x-axis (longest dimension), 275 and maintains periodic edges across the y- and z-boundaries. Each slice contains atoms and edges 276 within a fixed interval between x_0 and $x_0 + t_{slice}$, where x_0 and t_{slice} are the starting x-coordinate 277 and length of the slice, respectively. Atoms outside of a given slice but present in the connectivity 278 lists of those within are represented by virtual nodes (Fig. 4 - dashed circles). They are connected 279 to the slice using one-way virtual edges (Fig. 4 - dashed lines). Details about the construction of 280 partitions are given in Appendix C.

281 These virtual nodes/edges are initialized similarly to their 282 labeled counterparts with input atomic numbers and dis-283 tances. However, their outputs are masked and omitted 284 from the loss computation during training, validation, and 285 inference. Our masking approach differs from the one 286 used in transductive learning schemes, where the objec-287 tive is to learn the outputs of masked nodes/edges (Kipf & Welling (2016)). Here, we do not attempt to learn or 288 predict their outputs. The purpose of masked nodes and 289 edges here is to inform each partition of its full connectiv-290 ity and thus provide a much closer representation of the 291 graph topology. As the set of virtual connections used to 292 augment each graph corresponds only to the 1-hop neigh-293 borhood, we include only 1 MP layer in the network. During message passing, the network can then learn an 295 accurate and generalizable aggregation function when fit-296 ting to the output values of the labeled nodes and edges. 297 Hence, the network trained on a batch of such slices can predict the full graph of an unseen test structure. 298

299 300 301

3.3 A-HFO₂ STRUCTURE CREATION

To generate sufficiently rich training data, existing datasets typically sample molecules at various time steps of molecular dynamics (MD) trajectories (Yu et al. (2023a); Schütt et al. (2019); Christensen & lilienfeld (2020)) or generate multiple small perturbations of the atoms in a crystalline lattice (Li et al. (2022)). In the case of amorphous crystals, we take advantage of the fact that (1) almost each node has a different local atomic environ-



Figure 5: (a) Example atomic structure of a-HfO₂, showing oxygen atoms in blue and hafnium atoms in orange. The unit cell, which is illustrated by the black dashed box, has *x*-, *y*-, and *z*-dimensions of 53.876 Å × 26.308 Å × 26.242 Å, respectively. A slice (partition), characterized by a length t_{slice} , is also illustrated. (b) Distribution of the maximum element of each block of *H* as a function of interatomic distance, demonstrating the near-sighted nature of orbital interactions.

³⁰⁹ (1) annost each node has a different rotat atomic environ ³¹⁰ ment, and (2) the structure contains a large sampling of different motifs. A wide range of training ³¹⁰ data can thus be captured within a single sample. We, therefore, generated a dataset of only three ³¹¹ unique amorphous HfO₂ (a-HfO₂) structures for training, validation, and testing, respectively. The ³¹² DFT-based Hamiltonian of each of these systems (H^{GT}) was computed with a single- ζ valence ³¹³ (SZV) basis with 10 (4) Gaussian orbitals per Hf (O) atom. Details of the structure generation can ³¹⁴ be found in **Appendix F.1**. A structure example is shown in **Fig. 5**.

Structure	# atoms	# orbitals	# edges	x [Å]	y [Å]	z [Å]	nnz_H
1-validation	3,000	18,000	527,348	52.876	26.308	26.242	28,625,310
2-training 3-testing	3,000 3,000	18,000 18,000	533,364 530,920	52.346 52.722	26.237 26.267	26.293 26.191	28,943,862 28,805,422

318 319 320

315 316 317

Table 1: Attributes of the three generated a-HfO₂ structures: The [x, y, z] triplet defines the periodic unit cell size. nnz_H is the number of non-zero elements in the Hamiltonian, encompassing all orbital interactions. Edges were defined according to an interaction distance of $r_{cut} = 8$ Å.

³²⁴ 4 RESULTS 325

For fair comparisons in experiments where the quantity of training data varies, we used a ReduceL-RonPlateau scheduler that tracks the validation loss per epoch and reduces the learning rate by a fixed decay factor when no further decrease is detected. Training is stopped once a minimum learning rate is reached. Details on the values of the hyperparameters and the scheduler settings for different experiments can be found in **Appendix E**.

When trained on the $(H_{i,j})_{\alpha,\beta}$ elements of H₂O molecules from the MD17 dataset (Schütt et al. (2019)), our network as detailed in **Appendix A** can achieve prediction accuracy within an order of magnitude $(100 \times 10^{-6} E_h (E_h = \text{Hartree}) \text{ vs. } \sim 10 \times 10^{-6} E_h)$ of state-of-the-start equivariant GNN approaches, while using fewer layers (2 vs. 4-5) (**Appendix B**).

335 The treatment of the a-HfO₂ structures is, however, more challenging. We test the model's abil-336 ity to generalize to different configurations of this system by predicting H of the third structure in 337 **Table 1**, which remains unseen during the training process. The *augmented partitioning* scheme 338 is only applied during training (on structure 2), while the H of the unseen structure as a single 339 graph, including all nodes and edges, is predicted during inference (structure 3). Errors are reported 340 separately for nodes (ϵ_n) and edges (ϵ_e) for a more complete analysis that distinguishes intra- and 341 inter-atomic orbital interactions, which have very different magnitudes. Considering the large number of different motifs in the disordered structure, we also report the standard deviation of the node 342 and edge errors (σ_n and σ_e) to provide information on the consistency of the predictions. 343

344345 4.1 CUTOFF RADIUS AND CONNECTIVITY

346 We first explore the minimal graph connectivity that can be used by the network to accurately learn 347 relevant features. To do this we use the slice partition approach introduced in Section 3.2, using a 348 single slice of length $t_{slice} = 3$ Å to train the network. Reducing the value of r_{cut} below 8 Å no-349 ticeably increases the error ($\epsilon_{n/e}$), thus demonstrating the sensitivity of H to the exact connectivity 350 of the graph. Going from $r_{cut} = 8$ Å to 10 Å, the prediction error begins to plateau, but the node 351 degree (which is proportional to the memory consumption of the network) grows by $1.7 \times$. An r_{cut} 352 of 8 Å also results in negligible changes to the eigenvalue spectra (Fig. 10 in Appendix F.1). We 353 thus set r_{cut} =8 Å when defining graph edges for subsequent experiments. Note that $\epsilon_{node} >> \epsilon_{edge}$ 354 as the magnitude of the node labels is $\sim 100 \times$ larger than that of the edge labels. 355

r_{cut} [Å]	deg(n)	deg(n)'	Epochs	$\epsilon_{\mathbf{n}}$	$\sigma_{\mathbf{n}}$	$\epsilon_{\mathbf{e}}$	$\sigma_{\mathbf{e}}$
4	20.99	10.81	13,816	4.14	0.00960	9.60	0.00105
6	74.23	27.78	14,071	3.79	0.00925	0.40	0.00074
8	177.09	51.17	18,245	3.76	0.00903	0.22	0.00050
10	346.03	81.25	22,463	3.76	0.00886	0.15	0.00040

Table 2: Prediction accuracy of the network with different r_{cut} . Training was done with a single slice of length $t_{slice} = 3$ Å taken from structure 2 at $x_0 = 25$ Å. The edge connectivity of the matrix is set by r_{cut} . deg(n) is the average node degree, and deg(n)' the reduced node degree omitting virtual node neighbors. Note that for this value of t_{slice} , the majority of neighbors for the average node are virtual. ϵ_n and ϵ_e are the Mean Average Error (MAE) for nodes/edges, respectively, and σ_n and σ_e are the corresponding standard deviations. All units are in $[\times 10^{-3}E_h]$. In all cases, one MP layer is used. The validation loss of the model is computed from a slice of similar length, interaction distance, and starting location extracted from structure 1. The networks are tested on an unseen full graph (structure 3) constructed with the same r_{cut} .

369 370

356 357

359 360 361

4.2 ABLATION STUDIES OF THE TRAINING APPROACH

Next, we perform a study on the design features of the *augmented partitioning* approach introduced in **Section 3.2**. In particular, we examine the influence of virtual nodes, one-way/two-way, and one/two MP layers on the prediction accuracy. **Table 3** compiles ablation studies of the three aforementioned parameters, using 18 slices of length $t_{slice} = \sim 3$ Å that cover the full training structure (structure 2).

377 Compared to training with raw partitions, the addition of virtual nodes and edges reduces both ϵ_{node} and ϵ_{edge} by over ~50% when evaluated on the full test structure. Such an improvement is

378	Edges	# MP	$\epsilon_{\mathbf{n}}$	$\sigma_{\mathbf{n}}$	$\epsilon_{\mathbf{e}}$	$\sigma_{\mathbf{e}}$
380	n' n	1	5.18	7.93	1.66	0.0026
381	$n' \rightarrow n$	1	2.29	5.09 5.26	0.20	0.0038
382	$\frac{n \leftrightarrow n}{n' \rightarrow n}$	2	2.30	21.01	0.23	0.0042
383	$n \to n$ $n' \leftrightarrow n$	$\frac{2}{2}$	8.44	21.01 20.65	0.24	0.0043
384						

t_{slice} [Å]	N_t	N_e	Epochs	$\epsilon_{\mathbf{n}}$	$\sigma_{\mathbf{n}}$	$\epsilon_{\mathbf{e}}$	$\sigma_{\mathbf{e}}$
~ 1	54	47,958	15,744	2.30	5.37	0.20	0.37
~ 2	27	95,398	15,628	2.32	5.19	0.20	0.36
~ 3	18	141,512	15,675	2.29	5.09	0.20	0.38
~ 4	14	184,730	14,833	2.45	5.47	0.21	0.38
${\sim}8$	7	320,324	20,599	2.45	5.41	0.17	0.34
~ 12	5	381,504	19,351	2.59	5.86	0.18	0.36
\sim 52	1	533,364	23,396	2.46	5.39	0.16	0.33

385 Table 3: Ablation studies exploring the impact of 386 virtual nodes, one-way edges, and # of MP layers 387 on the prediction accuracy when tested on structure 3, using slices of length $t_{slice} = 3$ Å from 388 structure 2 for training and of length $t_{slice} = 4$ Å 389 (taken at $x_0 = 25$ Å) from structure 1 for vali-390 dation. The first column indicates the edge direc-391 tion between virtual (n') and labeled (n) nodes. 392 $n' \rightarrow n$ are one-way (incoming) edges, while $n^\prime \leftrightarrow n$ are two-way edges. Values are reported 393 in $\epsilon_n[mE_h]$ and $\sigma_n[\mu E_h]$. 394

Table 4: Prediction accuracy when the network is trained on differently-sized partitions of the same graph (structure 2), using 1 MP layer. $N_e = \#$ of slices, $N_e = \text{total } \#$ of labeled edges. The total number of labeled nodes remains constant. The number of slices is equal to $\sim L/t$, where L = 52.346 Å is the full length of structure 2 used for training. The validation set is a slice of length $t_{slice} =$ 4 Å starting at $x_0 = 25$ Å from structure 1. The models are tested on the full unseen structure (structure 3) and values are reported in $\epsilon_n [mE_h]$ and $\sigma_n [\mu E_h]$.

expected, as raw partitions are characterized by a large proportion of missing edges. We note that the best performing case achieves an ϵ_n and ϵ_e of 2.29 $\times 10^{-3}E_h$ and 0.20 $\times 10^{-3}E_h$, respectively, 397 398 using a single MP layer and one-way edges from virtual to labeled nodes $(n' \rightarrow n)$. As $t_{slice} < r_{cut}$, 399 a one-hop neighborhood is sufficient to cover all nodes across a partition. This property renders a 400 single MP layer sufficient, considering the near-sightedness of the Hamiltonian. The single MP-401 network's performance is trivially unaffected by the use of one-way virtual edges rather than two 402 $(n' \leftrightarrow n)$; virtual nodes do not need to aggregate information as they are omitted from the loss. 403 The edges connecting the labeled to virtual nodes $(n \rightarrow n')$ can thus be omitted to reduce memory 404 consumption. Using 2 MP layers with one layer of virtual nodes, however, degrades the performance 405 as the 2-hop neighborhood is not correctly represented. It should be emphasized that the full graph 406 with 2 MP layers does not fit into the memory of a single NVIDIA A100 GPU (requiring > 80 GiB). 407 Results with 2 MP layers are in **Appendix D** (under **Table 9**). 408

409 4.2.1 EFFECT OF AUGMENTED PARTITIONING ON PREDICTION ACCURACY

The *augmented partitioning* approach introduced in **Section 3.2** allows for the subdivision of the large graphs associated with the training of disordered materials by defining small slices (i.e., small t_{slice}) that can be treated sequentially on a single GPU or distributed across multiple GPUs and processed independently in parallel. Intermediate quantities do not need to be communicated during the forward/backward passes.

To establish that the network trained on augmented partitions does not suffer from loss of accu-416 racy, we partition the same full graph into different numbers of slices with different thicknesses for 417 training. In each case, the total number of labeled atoms summed up across all the slices remains 418 the same (3,000), while the total number of labeled edges reduces with increasing partitions. From 419 **Table 4** it can be seen that the prediction error is insensitive to partition size. Despite the different 420 divisions ranging from 5 ($t_{slice} = \sim 12$ Å) to 54 ($t_{slice} = \sim 1$ Å) slices, ϵ_n and ϵ_e remain very close to 421 the values obtained by training with the full graph ($t_{slice} = 52.346$ Å). For small slices, the reduced 422 fraction of labeled connections along the x direction does not affect the accuracy as the remaining 423 data along the y and z directions is sufficient to train the network. The combined MAE loss from 424 all 3,000 atoms and 530,920 edges for the best performing case ($t_{slice} = \sim 3$ Å) is 0.2159 mE_h, or 425 5.87 meV. This value is comparable to what a previous study obtained (2.2 meV) using equivariant GNNs for much smaller structures with \leq 150 atoms per unit cell (Wang et al. (2024b)). 426

427

429

395

410

428 4.3 PERFORMANCE ON A-HFO₂

To assess whether the prediction accuracy of the trained network is sufficient for practical application, we assemble the full Hamiltonian of the a-HfO₂ test structure using the network outputs (H^{pred}) , extract key quantities, and compare them with results obtained from the ground-truth



442 Figure 6: (a) ϵ_n of the predicted node blocks (red dashed lines) plotted against the full distribution of entries 443 in the ground-truth Hamiltonian matrix $(H_{i,j})_{\alpha,\beta}^{GT}$. (b) Same as (a), but for ϵ_e of the predicted edge blocks. (c) Eigenvalue spectra of the predicted (H^{pred}) and reference (H^{GT}) Hamiltonian matrices. The alpha value 444 indicates the scatter point transparency, such that the differences between the two scatter point sets can be more clearly seen. $(H_{i,j})^{pred}$ is symmetrized before diagonalization with $H = \frac{1}{2}(H + H^{\dagger})$. The relative 445 446 L1/L2 errors in the eigenvalue spectra, computed as $(\|(\vec{E})^{pred} - (\vec{E})^{GT}\|_{norm}^2)/(\|(\vec{E})^{GT}\|_{norm}^2)$ (norm = 447 1, 2) where \vec{E} is the vector of eigenvalues, are also shown for all eigenvalues and the ones corresponding 448 to occupied states (below 0.305 E_h). The black dashed box indicates the bandgap, which is defined as the 449 difference between the first energy above and below the Fermi level. 450

452 Hamiltonian (H^{GT}) computed with DFT. In Fig. 6 we show the achieved ϵ_n and ϵ_e relative to the magnitude of the H^{GT} entries. We then compute the eigenvalue spectrum of H^{pred} and H^{GT} 453 454 as well as the error distribution between them. The network used to create H^{pred} was trained with 455 18 partitions of length $t_{slice} = -3$ Å. We observe that both ϵ_n and ϵ_e are at least one order of magnitude lower than the average matrix element corresponding to the diagonal/off-diagonal blocks of the 456 Hamiltonian (black dashed line). This allows for the reconstructed H^{pred} to reproduce all eigen-457 values of H^{GT} within 0.87% relative L1 error. The error is 0.84% when eigenvalues of unoccupied 458 states above the cutoff 0.306 E_h are excluded. The remaining error is carried mostly by the largest 459 eigenvalues, and distributed around the edges of energy gaps (see Fig. 13 in Appendix H), which 460 correspond to regimes of stronger inter-atomic orbital coupling (Atkins & De Paula (2009)). 461

462 463

OTHER EXAMPLES 4.4

464 To further demonstrate the strength and robustness of the *augmented partitioning approach*, we 465 have conducted additional experiments on other datasets (found in Appendix I). There we showed 466 that the approach generalizes well to HfO_x structures with differently distributed vacancies, and the 467 model achieves an even better prediction accuracy of 1.43 meV for amorphous PtGe, compared to 468 the 5.87 meV achieved for HfO₂. 469

- 4.5 COMPUTATIONAL COST 470
- 471

Compared to a naive full-batch training of the graph, our method using just 8 augmented slices 472 results in a $6.5 \times$ speedup per epoch (0.38 vs. 2.5 s), and a $7.2 \times$ decrease in memory consumption 473 per rank (8.59 vs. 61.68 GiB). A more complete analysis is provided in Appendix G.2. This 474 scaling behavior is limited only by the overhead introduced by the virtual nodes/edges, and the 475 small computational load imbalance from partitioning. Further computational improvements could 476 be achieved by combining the augmentation approach with optimized graph partitioning algorithms 477

(Karypis & Kumar (1998)) while leveraging periodicity.

478 The extension of GNN-based predictions to large material systems could potentially save tremen-479 dous amounts of computational time. While DFT calculations to obtain the H^{GT} of small molecules 480 (e.g., H₂O) take only a few seconds, the same operation for a-HfO₂ structures made of 3,000 atoms 481 is computationally $\sim 100 \times$ heavier (~ 0.04 vs. ~ 3.65 node hours, see Appendix G). More im-482 portantly, the GNN prediction unlocks the ability to consider much larger structures than the ones 483 considered here, the inference phase scaling with $O(N_{atoms})$ while DFT calculations are limited to 484 $\mathcal{O}(N_{atoms}^3)$. The model could also serve as an initial guess to DFT packages to reduce the number 485 of self-consistent field iterations that are required to obtain converged electron densities (Unke et al. (2021)).

486 5 CONCLUSION

488 We adapted equivariant GNNs to learn the electronic properties of amorphous materials, and intro-489 duced an *augmented partitioning* approach to break down and train the large graphs encountered 490 when dealing with realistic structural disorder, without sacrificing accuracy. More generally, we 491 proposed a method to tackle the training of atomic systems that require large, highly connected, and 492 near-sighted graphs where a strictly local atomic environment is sufficient. The key is the addition of 493 virtual nodes and edges connected to relatively small partitions that mimic their neighborhood. The 494 method, demonstrated here on a-HfO₂, can be straightforwardly applied to other disordered materials, or adapted to learn their other rotationally-equivariant attributes such as vibrational properties, 495 e.g., phonon dispersions (Fang et al. (2024)). 496

497 The resulting networks capture relevant properties of a-HfO₂ in sufficient detail to achieve few-meV498 accuracy and reproduce features of practical relevance, such as the energy eigenvalues. However, 499 the sub-millielectronvolt range is not currently reached, contrary to what has been demonstrated 500 with small-molecule datasets (Yu et al. (2023b); Unke et al. (2021)). This shortcoming can most likely be attributed to a combination of greater dataset complexity in amorphous compounds, coarser 501 resolution of the training data ($\epsilon_{SCF} = 1 \times 10^{-6} E_h$), and limited network size. Note that the 502 augmented partitioning approach is also a general method that can also be adapted for use in other, 503 more expressive network architectures, since it is mainly applied during graph construction. Further 504 data generation, parameter optimization, and enabling of networks with increased expressiveness 505 will be the next steps.

506 507 508

5.1 OUTLOOK & APPLICATIONS

509 The ability to learn the electronic properties of complex disordered materials unlocks notable appli-510 cations in computational physics, chemistry, and materials science. Several compounds are used in 511 their amorphous phase, as they often exhibit different properties from their crystalline equivalents. 512 For example, $a-SiO_2$ as dielectric layer has been a key enabler of the metal-oxide-semiconductor 513 technology (Nekrashevich & Gritsenko (2014)), IGZO, thanks to its large electron mobility, serves 514 as channel of flexible transistors (Kamiya et al. (2010)), a-HfO₂ allows for the (non-)volatile storage 515 of information, when placed between two metallic electrodes (Chan et al. (2008)), and GST can be used to write/store data optically (Pirovano et al. (2004); Kolobov et al. (2004)). Downscaling of 516 materials also reveals structural effects similar to disorder, e.g., defects (Wilhelmer et al. (2022)), 517 strain Parton & Verheyen (2006), or grain boundaries (Weitz et al. (2009)). All of them require large 518 unit cells to be accurately described (Lany & Zunger (2008); Zhao et al. (2020)). Computationally 519 expensive DFT calculations represent a bottleneck towards investigating the electronic properties of 520 such systems in simulations. Recent advances in graph neural networks, combined with domain-521 specific innovations to train them, will enable such explorations. 522

- 523 524 AUTHOR CONTRIBUTIONS
- 525 [anonymized] 526
- 527 ACKNOWLEDGMENTS
- 529 [anonymized]
- 530

528

531 REPRODUCIBILITY 532

The code used to set up and train the network is available at the following repository: [anonymized code temporarily provided as supplementary material]. CP2K was used to generate the custom HfO₂ dataset. Input files are provided in the same repository.

536

537 REFERENCES

539 Vladimir I. Anisimov, Jan Zaanen, and Ole K. Andersen. Band theory and mott insulators: Hubbard u instead of stoner i. *Physical Review B*, 44:943–954, Jul 1991.

556

558

568

- Peter Atkins and Julio De Paula. *Atkins' physical chemistry*. Oxford University Press, London, England, 9 edition, November 2009.
- Hexin Bai, Peng Chu, Jeng-Yuan Tsai, Nathan Wilson, Xiaofeng Qian, Qimin Yan, and Haibin
 Ling. Graph neural network for hamiltonian-based material property prediction. *Neural Com- puting and Applications*, 34(6):4625–4632, November 2021. ISSN 1433-3058. doi: 10.1007/
 s00521-021-06616-0. URL http://dx.doi.org/10.1007/s00521-021-06616-0.
- Ilyes Batatia, Lars L. Schaaf, Huajie Chen, Gábor Csányi, Christoph Ortner, and Felix A. Faber.
 Equivariant matrix function neural networks, 2023. URL https://arxiv.org/abs/ 2310.10434.
- Maciej Besta and Torsten Hoefler. Parallel and distributed graph neural networks: An in-depth concurrency analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(5): 2584–2606, May 2024. ISSN 1939-3539. doi: 10.1109/tpami.2023.3303431. URL http://dx.doi.org/10.1109/TPAMI.2023.3303431.
 - M.Y. Chan, T. Zhang, V. Ho, and P.S. Lee. Resistive switching effects of hfo2 high-k dielectric. *Microelectronic Engineering*, 85(12):2420–2424, December 2008. ISSN 0167-9317. doi: 10. 1016/j.mee.2008.09.021. URL http://dx.doi.org/10.1016/j.mee.2008.09.021.
- Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery amp; Data Mining*, KDD '19. ACM, July 2019. doi: 10.1145/3292500.3330925. URL http://dx.doi.org/10.1145/3292500.3330925.
- J.H. Choi, Y. Mao, and J.P. Chang. Development of hafnium based high-k materials—a review. *Materials Science and Engineering: R: Reports*, 72(6):97–136, July 2011. ISSN 0927-796X. doi: 10.1016/j.mser.2010.12.001. URL http://dx.doi.org/10.1016/j.mser.2010.
 12.001.
- Anders S. Christensen and Anatole Von lilienfeld. Revised MD17 dataset (rMD17). 7 2020. doi: 10.
 6084/m9.figshare.12672038.v3. URL https://figshare.com/articles/dataset/
 Revised_MD17_dataset_rMD17_/12672038.
- Ralf Drautz. Atomic cluster expansion of scalar, vectorial, and tensorial properties including magnetism and charge transfer. *Physical Review B*, 102(2), July 2020. ISSN 2469-9969. doi: 10.1103/physrevb.102.024104. URL http://dx.doi.org/10.1103/PhysRevB.102.024104.
- Fabian Ducry, Jan Aeschlimann, and Mathieu Luisier. Electro-thermal transport in disordered nanostructures: a modeling perspective. *Nanoscale Advances*, 2(7):2648–2667, 2020. ISSN 2516-0230.
 doi: 10.1039/d0na00168f. URL http://dx.doi.org/10.1039/d0na00168f.
- Shiang Fang, Mario Geiger, Joseph G. Checkelsky, and Tess Smidt. Phonon predictions with e(3)-equivariant graph neural networks, 2024. URL https://arxiv.org/abs/2403.11347.
- Xiaoxun Gong, He Li, Nianlong Zou, Runzhang Xu, Wenhui Duan, and Yong Xu. General frame work for e(3)-equivariant neural network representation of density functional theory hamiltonian.
 Nature Communications, 14(1), May 2023. ISSN 2041-1723. doi: 10.1038/s41467-023-38468-8.
 URL http://dx.doi.org/10.1038/s41467-023-38468-8.
- Qiangqiang Gu, Zhanghao Zhouyin, Shishir Kumar Pandey, Peng Zhang, Linfeng Zhang, and Weinan E. Deep learning tight-binding approach for large-scale electronic simulations at finite temperatures with ab initio accuracy. *Nature Communications*, 15(1), August 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-51006-4. URL http://dx.doi.org/10.1038/ s41467-024-51006-4.
- Ganesh Hegde and R. Chris Bowen. Machine-learned approximations to density functional theory hamiltonians. *Scientific Reports*, 7(1), February 2017. ISSN 2045-2322. doi: 10.1038/srep42669. URL http://dx.doi.org/10.1038/srep42669.

594 P. Hohenberg and W. Kohn. Inhomogeneous electron gas. *Physical Review*, 136(3B):B864–B871, 595 November 1964. ISSN 0031-899X. doi: 10.1103/physrev.136.b864. URL http://dx.doi. 596 org/10.1103/PhysRev.136.B864. 597 Toshio Kamiya, Kenji Nomura, and Hideo Hosono. Present status of amorphous in-ga-zn-o thin-598 film transistors. Science and Technology of Advanced Materials, 11(4):044305, February 2010. ISSN 1878-5514. doi: 10.1088/1468-6996/11/4/044305. URL http://dx.doi.org/10. 600 1088/1468-6996/11/4/044305. 601 602 Manasa Kaniselvan, Mathieu Luisier, and Marko Mladenović. An atomistic model of field-induced 603 resistive switching in valence change memory. ACS Nano, 17(9):8281-8292, March 2023. ISSN 1936-086X. doi: 10.1021/acsnano.2c12575. URL http://dx.doi.org/10.1021/ 604 acsnano.2c12575. 605 George Karypis and Vipin Kumar. A parallel algorithm for multilevel graph partitioning and sparse 607 matrix ordering. Journal of Parallel and Distributed Computing, 48(1):71–95, January 1998. 608 ISSN 0743-7315. doi: 10.1006/jpdc.1997.1403. URL http://dx.doi.org/10.1006/ 609 jpdc.1997.1403. 610 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional net-611 works. CoRR, abs/1609.02907, 2016. URL http://arxiv.org/abs/1609.02907. 612 613 Cedric Klinkert, Aron Szabó, Christian Stieger, Davide Campi, Nicola Marzari, and Mathieu Luisier. 614 2-d materials for ultrascaled field-effect transistors: One hundred candidates under the ab initio 615 microscope. ACS Nano, 14(7):8605-8615, June 2020. ISSN 1936-086X. doi: 10.1021/acsnano. 616 0c02983. URL http://dx.doi.org/10.1021/acsnano.0c02983. 617 W. Kohn and L. J. Sham. Self-consistent equations including exchange and correlation effects. 618 Physical Review, 140(4A):A1133–A1138, November 1965. ISSN 0031-899X. doi: 10.1103/ 619 physrev.140.a1133. URL http://dx.doi.org/10.1103/PhysRev.140.A1133. 620 621 Alexander V. Kolobov, Paul Fons, Anatoly I. Frenkel, Alexei L. Ankudinov, Junji Tominaga, and 622 Tomoya Uruga. Understanding the phase-change mechanism of rewritable optical media. Nature 623 Materials, 3(10):703-708, September 2004. ISSN 1476-4660. doi: 10.1038/nmat1215. URL 624 http://dx.doi.org/10.1038/nmat1215. 625 Thomas D. Kühne, Marcella Iannuzzi, Mauro Del Ben, Vladimir V. Rybkin, Patrick Seewald, Fred-626 erick Stein, Teodoro Laino, Rustam Z. Khaliullin, Ole Schütt, Florian Schiffmann, Dorothea 627 Golze, Jan Wilhelm, Sergey Chulkov, Mohammad Hossein Bani-Hashemian, Valéry Weber, Ur-628 ban Borštnik, Mathieu Taillefumier, Alice Shoshana Jakobovits, Alfio Lazzaro, Hans Pabst, and et 629 al. CP2k: An electronic structure and molecular dynamics software package - quickstep: Efficient 630 and accurate electronic structure calculations. J. Chem. Phys., 152(19):194103, May 2020. 631 Stephan Lany and Alex Zunger. Assessment of correction methods for the band-gap problem and 632 for finite-size effects in supercell defect calculations: Case studies for zno and gaas. Physical 633 Review B, 78(23), December 2008. ISSN 1550-235X. doi: 10.1103/physrevb.78.235104. URL 634 http://dx.doi.org/10.1103/PhysRevB.78.235104. 635 636 He Li, Zun Wang, Nianlong Zou, Meng Ye, Runzhang Xu, Xiaoxun Gong, Wenhui Duan, 637 and Yong Xu. Deep-learning density functional theory hamiltonian for efficient ab initio 638 electronic-structure calculation. Nature Computational Science, 2(6):367-377, June 2022. ISSN 639 2662-8457. doi: 10.1038/s43588-022-00265-6. URL http://dx.doi.org/10.1038/ s43588-022-00265-6. 640 641 Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, 642 Jeff Smith, Brian Vaughan, Pritam Damania, and Soumith Chintala. Pytorch distributed: Experi-643 ences on accelerating data parallel training, 2020. URL https://arxiv.org/abs/2006. 644 15704. 645 Renjie Liao, Marc Brockschmidt, Daniel Tarlow, Alexander L. Gaunt, Raquel Urtasun, and Richard 646 Zemel. Graph partition neural networks for semi-supervised classification, 2018. URL https: 647

//arxiv.org/abs/1803.06272.

655

656

657

658

671

672

673

674

675

676

677 678

679

680

681

682

683

684

685

686

687

688 689

690

691

692

- Yi-Lun Liao, Brandon Wood, Abhishek Das, and Tess Smidt. Equiformerv2: Improved equivariant transformer for scaling to higher-degree representations, 2023. URL https://arxiv.org/abs/2306.12059.
- Frank Neese. The orca program system. WIREs Computational Molecular Science, 2(1):73–78, June 2011. ISSN 1759-0884. doi: 10.1002/wcms.81. URL http://dx.doi.org/10.1002/wcms.81.
 - S. S. Nekrashevich and V. A. Gritsenko. Electronic structure of silicon dioxide (a review). *Physics of the Solid State*, 56(2):207–222, February 2014. ISSN 1090-6460. doi: 10.1134/s106378341402022x. URL http://dx.doi.org/10.1134/s106378341402022x.
- Jigyasa Nigam, Michael J. Willatt, and Michele Ceriotti. Equivariant representations for molecular hamiltonians and n-center atomic-scale properties. *The Journal of Chemical Physics*, 156(1), January 2022. ISSN 1089-7690. doi: 10.1063/5.0072784. URL http://dx.doi.org/10. 1063/5.0072784.
- Els Parton and Peter Verheyen. Strained silicon the key to sub-45 nm cmos. *III-Vs Review*, 19 (3):28–31, April 2006. ISSN 0961-1290. doi: 10.1016/s0961-1290(06)71590-3. URL http://dx.doi.org/10.1016/S0961-1290(06)71590-3.
- Saro Passaro and C. Lawrence Zitnick. Reducing so(3) convolutions to so(2) for efficient equivariant gnns, 2023. URL https://arxiv.org/abs/2302.03655.
- John P. Perdew, Kieron Burke, and Matthias Ernzerhof. Generalized gradient approximation made simple [phys. rev. lett. 77, 3865 (1996)]. *Phys. Rev. Lett.*, 78:1396–1396, Feb 1997.
 - A. Pirovano, A.L. Lacaita, A. Benvenuti, F. Pellizzer, and R. Bez. Electronic switching in phasechange memories. *IEEE Transactions on Electron Devices*, 51(3):452–459, March 2004. ISSN 0018-9383. doi: 10.1109/ted.2003.823243. URL http://dx.doi.org/10.1109/TED. 2003.823243.
 - Chendi Qian, Andrei Manolache, Christopher Morris, and Mathias Niepert. Probabilistic graph rewiring via virtual nodes, 2024. URL https://arxiv.org/abs/2405.17311.
 - Gil M. Repa and Lisa A. Fredin. Predicting electronic structure of realistic amorphous surfaces. *Advanced Theory and Simulations*, 6(11), June 2023. ISSN 2513-0390. doi: 10.1002/adts. 202300292. URL http://dx.doi.org/10.1002/adts.202300292.
 - K. T. Schütt, M. Gastegger, A. Tkatchenko, K.-R. Müller, and R. J. Maurer. Unifying machine learning and quantum chemistry with a deep neural network for molecular wavefunctions. *Nature Communications*, 10(1), November 2019. ISSN 2041-1723. doi: 10.1038/s41467-019-12875-2. URL http://dx.doi.org/10.1038/s41467-019-12875-2.
 - M. L. Senent and S. Wilson. Intramolecular basis set superposition errors. *International Journal of Quantum Chemistry*, 82(6):282–292, 2001. doi: https://doi.org/10.1002/qua.1030. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/qua.1030.
 - Jack Strand, Moloud Kaviani, David Gao, Al-Moatasem El-Sayed, Valeri V Afanas'ev, and Alexander L Shluger. Intrinsic charge trapping in amorphous oxide films: status and challenges. *Journal* of Physics: Condensed Matter, 30(23):233001, May 2018. ISSN 1361-648X. doi: 10.1088/ 1361-648x/aac005. URL http://dx.doi.org/10.1088/1361-648X/aac005.
- Anders Blom Troels Markussen Jess Wellendorff Julian Schneider Tue Gunst Brecht Verstichel
 Petr A Khomyakov Ulrik G Vej-Hansen Mads Brandbyge Søren Smidstrup, Kurt Stokbro et al.
 Quantumatk: An integrated platform of electronic and atomic-scale modelling tools. J. Phys:
 Condens. Matter (APS), 32:015901, 2020. doi: 10.1088/1361-648X/ab4007. URL https:
 //iopscience.iop.org/article/10.1088/1361-648X/ab4007.
- De Nyago Tafen and D. A. Drabold. Realistic models of binary glasses from models of tetrahedral amorphous semiconductors. *Physical Review B*, 68(16), October 2003. ISSN 1095-3795.
 doi: 10.1103/physrevb.68.165208. URL http://dx.doi.org/10.1103/PhysRevB. 68.165208.

702 703 704	Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: Rotation- and translation-equivariant neural networks for 3d point clouds, 2018. URL https://arxiv.org/abs/1802.08219.
705 706 707 708	Oliver T. Unke, Mihail Bogojeski, Michael Gastegger, Mario Geiger, Tess Smidt, and Klaus-Robert Müller. Se(3)-equivariant prediction of molecular wavefunctions and electronic densities, 2021. URL https://arxiv.org/abs/2106.02347.
709 710 711 712 713	M. Laura Urquiza, Md Mahbubul Islam, Adri C. T. van Duin, Xavier Cartoixà, and Alejandro Stra- chan. Atomistic insights on the full operation cycle of a hfo2-based resistive random access memory cell from molecular dynamics. <i>ACS Nano</i> , 15(8):12945–12954, July 2021. ISSN 1936- 086X. doi: 10.1021/acsnano.1c01466. URL http://dx.doi.org/10.1021/acsnano. 1c01466.
714 715 716 717	Joost VandeVondele and Jürg Hutter. Gaussian basis sets for accurate calculations on molecular systems in gas and condensed phases. <i>J. Chem. Phys.</i> , 127(11):114105, 09 2007. ISSN 0021-9606.
718 719	Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. In <i>ICLR</i> , 2018.
720 721 722 723	Cheng Wan, Youjie Li, Cameron R Wolfe, Anastasios Kyrillidis, Nam Sung Kim, and Yingyan Lin. PipeGCN: Efficient full-graph training of graph convolutional networks with pipelined feature communication. March 2022.
724 725 726 727	Yuxiang Wang, He Li, Zechen Tang, Honggeng Tao, Yanzhen Wang, Zilong Yuan, Zezhou Chen, Wenhui Duan, and Yong Xu. Deeph-2: Enhancing deep-learning electronic structure via an equivariant local-coordinate transformer, 2024a. URL https://arxiv.org/abs/2401. 17015.
728 729 730 731 732 733	Yuxiang Wang, Yang Li, Zechen Tang, He Li, Zilong Yuan, Honggeng Tao, Nianlong Zou, Ting Bao, Xinghao Liang, Zezhou Chen, Shanghua Xu, Ce Bian, Zhiming Xu, Chong Wang, Chen Si, Wenhui Duan, and Yong Xu. Universal materials model of deep-learning density functional theory hamiltonian. <i>Science Bulletin</i> , 69(16):2514–2521, 2024b. ISSN 2095-9273. doi: https://doi.org/10.1016/j.scib.2024.06.011. URL https://www.sciencedirect.com/science/article/pii/S2095927324004079.
734 735 736 737	Yuyang Wang, Zijie Li, and Amir Barati Farimani. Graph Neural Networks for Molecules, pp. 21–66. Springer International Publishing, 2023. ISBN 9783031371967. doi: 10.1007/978-3-031-37196-7_2. URL http://dx.doi.org/10.1007/978-3-031-37196-7_2.
738 739 740 741 742 743	R. T. Weitz, K. Amsharov, U. Zschieschang, M. Burghard, M. Jansen, M. Kelsch, B. Rhamati, P. A. van Aken, K. Kern, and H. Klauk. The importance of grain boundaries for the time-dependent mobility degradation in organic thin-film transistors. <i>Chemistry of Materials</i> , 21(20):4949–4954, September 2009. ISSN 1520-5002. doi: 10.1021/cm902145x. URL http://dx.doi.org/10.1021/cm902145x.
744 745 746 747 748	Christoph Wilhelmer, Dominic Waldhoer, Markus Jech, Al-Moatasem Bellah El-Sayed, Lukas Cvitkovich, Michael Waltl, and Tibor Grasser. Ab initio investigations in amorphous silicon dioxide: Proposing a multi-state defect model for electron and hole capture. <i>Microelectronics Reliability</i> , 139:114801, December 2022. ISSN 0026-2714. doi: 10.1016/j.microrel.2022.114801. URL http://dx.doi.org/10.1016/j.microrel.2022.114801.
749 750 751 752 753	Yong Youn, Youngho Kang, and Seungwu Han. An efficient method to generate amorphous struc- tures based on local geometry. <i>Computational Materials Science</i> , 95:256–262, December 2014. ISSN 0927-0256. doi: 10.1016/j.commatsci.2014.07.053. URL http://dx.doi.org/10. 1016/j.commatsci.2014.07.053.
754	Haiyang Yu, Meng Liu, Youzhi Luo, Alex Strasser, Xiaofeng Qian, Xiaoning Qian, and Shuiwang Ji Obly A quantum hamiltonian prediction benchmark for gm9 molecules 2023a URL https://

Ji. Qh9: A quantum hamiltonian prediction benchmark for qm9 molecules, 2023a. URL https: //arxiv.org/abs/2306.09549.

756 757 758	Haiyang Yu, Zhao Xu, Xiaofeng Qian, Xiaoning Qian, and Shuiwang Ji. Efficient and equivariant graph networks for predicting quantum hamiltonian, 2023b. URL https://arxiv.org/abs/2306.04922.
759 760 761 762	He Zhang, Chang Liu, Zun Wang, Xinran Wei, Siyuan Liu, Nanning Zheng, Bin Shao, and Tie-Yan Liu. Self-consistency training for density-functional-theory hamiltonian prediction, 2024. URL https://arxiv.org/abs/2403.09560.
763 764 765 766 767	Liwei Zhang, Berk Onat, Geneviève Dusson, Adam McSloy, G. Anand, Reinhard J. Maurer, Christoph Ortner, and James R. Kermode. Equivariant analytical mapping of first principles hamiltonians to accurate and transferable materials models. <i>npj Computational Materials</i> , 8(1), July 2022. ISSN 2057-3960. doi: 10.1038/s41524-022-00843-2. URL http://dx.doi.org/10.1038/s41524-022-00843-2.
768 769 770 771	Xin-Gang Zhao, Gustavo M. Dalpian, Zhi Wang, and Alex Zunger. Polymorphous nature of cubic halide perovskites. <i>Physical Review B</i> , 101(15), April 2020. ISSN 2469-9969. doi: 10.1103/ physrevb.101.155137. URL http://dx.doi.org/10.1103/PhysRevB.101.155137.
772 773 774 775	Yang Zhong, Hongyu Yu, Mao Su, Xingao Gong, and Hongjun Xiang. Transferable equivariant graph neural networks for the hamiltonians of molecules and solids. <i>npj Computational Materials</i> , 9(1), October 2023. ISSN 2057-3960. doi: 10.1038/s41524-023-01130-4. URL http://dx. doi.org/10.1038/s41524-023-01130-4.
776 777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
000 807	
808	
809	

810 A NETWORK ARCHITECTURE



831 832

838

812

Figure 7: The node embeddings after the message passing layer labeled X are denoted as n_X , and the corresponding edge features are e_X . a_{ij} is the block of attention weights, Z_i is the atomic number of atom *i*, and r_{ij} is the scalar distance between atoms *i* and *j*. The input message to the update blocks consist of the concatenated embeddings of source and target nodes *j* and *i*, along with their connecting edge e_{ji} . They are multiplied by a set of weights generated from the radial function using a scalar embedding that consists of the atomic numbers Z_i and Z_j concatenated with the edge distance $|r_{ji}|$

In this section, we provide further details on the architecture adapted from EquiformerV2 (Liao et al. (2023)) and DeepH2 (Wang et al. (2024a)), in addition to the network initialization description in Section 3.1.

B42 During message passing, input messages in the form of node embeddings of size $(N_n, (l_{max} + 1)^2, d_{sphere,n})$ and edge embeddings of size $(N_e, (l_{max} + 1)^2, d_{sphere,n})$ from surrounding neighboring B44 source nodes are passed into the target node. Note that for our network, $d_{sphere,n}$ and $d_{sphere,e}$ are B45 set to the same value d_{sphere} , known as the number of spherical channels.

Inside the node update block, the source and target atom embeddings are concatenated together with the edge embeddings of their connecting edge to form an input message of size $(N_e, (l_{max} + 1)^2, 3d_{sphere})$. The message is multiplied by weights generated by the radial functions using scalar embeddings (atomic numbers and distances), making it more receptive to small changes in environments. The dimension of these scalar embeddings are similarly set to d_{sphere} .

851 Afterwards, the input message is rotated to align with the z axis through the rotation block. The ro-852 tation block contains rotation matrices that were pre-computed using the normalized edge vectors of 853 every edge in the graph. The rotated message is reshaped into m major order for linear convolutions 854 to be performed for each m = 0 to $m = m_{max}$ with the number of l components for each m is given by $l_{max} - m + 1$. m_{max} and l_{max} ($m_{max} \le l_{max}$) are other hyperparameters that can be adjusted. 855 The convolutions produce an output embedding of size $(n_e, (l_{max} + 1)^2, d_{attn.hidden})$, where atten-856 tion hidden channels $d_{attn.hidden}$ is another hyperparameter. The message is then fed through the gate activation layer, which adds non-linearity while preserving equivariance by applying separate 858 non-linearities to the l = 0 and l > 0 components. 859

Next, the non-linear message is passed through a second convolution layer that produces an output embedding size of $(n_e, (l_{max}+1)^2, d_{attn_value}*N_{heads})$, which is then be reshaped into $(n_e, (l_{max}+1)^2, N_{heads}, d_{attn_value})$. For each edge surrounding the target node, a set of d_{attn_alpha} attention weights are then generated for each attention head, with the total number of heads being N_{heads} . This is used to generate the output vector alpha, which is reshaped along N_{heads} and multiplied with the reshaped message embedding. Finally, the messages from neighboring are reshaped, rotated back, and aggregated onto the target node, before being projected back into the shape $(N_n, (l_{max} + 1)^2, d_{sphere})$.

This is finally passed into a feed-forward network consisting of two linear layers that mixes features of the same l together. The hidden dimensions used in the feed-forward network is given by the hyperparameter d_{ffn} . The edge update block is similar to the node update block, except there is only one convolution layer and no attention required. Between the blocks, layer normalisation is also applied, and similar to EquiformerV2, we also normalised the l = 0 features separately from the l > 0 features. The final predicted output is passed into the Wigner Eckart layer to be reconstructed into Hamiltonian blocks.

874

875 876 A.1 WIGNER ECKART LAYER

877 Our implementation of this layer is similar to that found in Gong et al. (2023). A Hamiltonian block 878 representing the interaction between atom *i* and *j* consists can be split into different sub-blocks 879 representing interactions between orbital 1 of order l_1 and orbital 2 of degree l_2 as shown in Fig. 1. 880 Each block of size $(2l_1 + 1) \times (2l_2 + 1)$, is an equivariant tensor that comes from the tensor product 881 $l_1 \otimes l_2$ for every pair of interacting orbitals in atom *i* and atom *j*. For example, an interaction block 882 between a *p* orbital (l = 1) and a *d* orbital (l = 2) has $(2 \times 1 + 1) \times (2 \times 2 + 1) = 3 \times 5$ elements. 883 We refer to this form as the uncoupled basis representation of the Hamiltonian block.

Before training, the Wigner-Eckart layer converts the Hamiltonian data from uncoupled basis to
 coupled basis using Clebsch Gordan coefficients.

886
$$l_1 \otimes l_2 = |l_1 - l_2| \oplus ... \oplus (l_1 + l_2)$$

888 $l_1 \otimes l_2$ is now represented by a direct sum of coupled sub-spaces with order ranging from $|l_1 - l_2|$ to 889 $(l_1 + l_2)$. This is repeated for every possible type of orbital interaction between the atoms, to obtain 890 the final direct sum of all the sub-spaces needed for the network to reconstruct the Hamiltonian.

The largest possible $(l_1 + l_2)$ value determines the minimum l_{max} hyperparameter needed for the embeddings of the equivariant model. In the case of HfO₂ in this paper, which uses the SZV basis set, the highest order orbital is the *d* orbital of Hf. This means that the largest possible (l_1+l_2) comes from the interaction between two Hf *d* orbitals, and the l_{max} needs to be at least equal to (2+2) = 4. This allows the predicted output of the model to be converted back into the full uncoupled basis using the same layer, and reassembled into the full Hamiltonian matrix during inference.

897

899

A.2 LOSS COMPUTATION DURING TRAINING AND INFERENCE

For all experiments, a minor difference from the procedures reported in Yu et al. (2023b) and Schütt et al. (2019) is that we use the Mean Squared Error (MSE) of the full target vectors in the coupled space to compute the fitting and validation loss during training:

$$\mathcal{L}_{MSE}(x_i, \hat{x}_i) = \frac{1}{N} \sum_{i=1}^{N} (x_i - \hat{x}_i)^2$$

Where x and \hat{x} are the flattened targets (orbital blocks). These targets are padded with zeros to ensure that those of different orbital interactions have the same dimensions. To avoid re-shaping the predictions within every epoch, the loss computed during training includes this padding. This is procedure was also used in DeepH E3 (Gong et al. (2023)). However, the final reported loss in Table 6 uses the Mean Absolute Error (MAE) after converting the output and label tensors back into uncoupled space and reconstructing the Hamiltonian blocks.

$$\mathcal{L}_{MAE}(x_i, \hat{x}_i) = \frac{1}{N} \sum_{i=1}^{N} |x_i - \hat{x}_i|$$

915 916

913 914

917 The padding is omitted during the reconstruction process, and all the elements in the blocks where the final MAE is computed from represent orbital interactions that exist in the label Hamiltonian.

929

930

931 932 933

934 935

936

937

938

939

940

941

942

943 944

946

947

957

958

959

960



Figure 8: Training the MD17 Hamiltonian matrices of H₂O molecules. (a) Elements of the hamiltonian matrix for a sample H_2O molecule in the dataset. (b) Fitting of the predicted elements of the Hamiltonian matrix against the labels, after randomized downsampling to every 100 data points. Data points closer to the grey line indicate better agreement. (c) Magnification of fit around [-1, 1].

В PERFORMANCE ON SMALL MOLECULES (H_2O)

As our work combines and adapts existing methodology to a new, application-oriented dataset, we first ensure that our adaptation of the networks introduced by Liao et al. (2023) and Wang et al. (2024a) can achieve reasonable results on an existing dataset. We select the dataset of H_2O molecules in MD17 Schütt et al. (2019). The hyperparameters which determine the network size are d_{sphere} and d_{attn_hidden} (Table 8). In Table 5 we explore the effect of these parameters under early-stopping conditions. We fix the number of MP layers to 2 in all cases, lower than was used by similar equivariant GNNs (Unke et al. (2021), Yu et al. (2023b)). Channel dimensions of 32 and above meet the loss defined for the stopping criteria, with larger dimensions reaching it in fewer epochs.

We select a dimension of 64, and train the network under the conditions specified in **Table 8**. In **Ta**-945 **ble 6** we compare the final Mean Absolute Error (MAE) loss with similar networks in the literature. Our implementation can achieve a prediction accuracy of $\sim 100 \ \mu E_h$, which is within an order of magnitude of other equivariant networks while using a smaller network (#MP layers). 948

Dim.	# Epochs	Training I	Molecules	Testing M	Iolecules
		$\epsilon_{node} \left[\times 10^{-6} E_h \right]$	$\epsilon_{edge} \left[\times 10^{-6} E_h \right]$	$\epsilon_{node} \left[\times 10^{-3} E_h \right]$	$\epsilon_{edge} \left[\times 10^{-3} E_h \right]$
16/16	50,000*	0.20*	2.81*	0.1498*	0.4776*
32/32	10,847	0.54	0.43	0.1716	0.1866
64/64	6,761	0.36	0.32	0.1322	0.1534
28/128	2,015*	1.08*	1.01*	0.1878*	0.2776*

Table 5: Effect of the channel dimension used to train the network on H₂O. For this study we remove 1499 molecules from the dataset (4999 total molecules), and use 500 for training (single-batch), 500 for validation, and 499 for testing. For fair comparison between the different network sizes, the trainings were subject to an early-stopping criteria of $\epsilon_{MSE} = 1 \times 10^{-6}$ of the validation set loss. Training losses are reported in MSE while testing losses are reported in MAE. The '*' indicates that the error on the validation set did not reach the stopping criteria, and training instead finished when $lr = \epsilon_{MSE}$ (128/128) or when the # epochs reached 50,000.

Network	MAE H $[meV]$	MAE H [$\times 10^{-6} E_h$]
PhiSNet	0.47	18
QHNet	0.29	11
This work	2.7	100

968 Table 6: Mean Absolute Error (MAE) for predicted Hamiltonian matrices reported for PhiSNet (Unke et al. 969 (2021)), and QHNet (Yu et al. (2023b)) on the MD17 H₂O dataset, taken from the respective publications. The 970 result for our adapted EquiformerV2 network for electronic property prediction uses 500 molecules for testing, 500 for validation, 2500 for testing. Note that we use 2 layers, as opposed to PhiSNet (4) and QHnet (5). 971

972 C AUGMENTING GRAPH PARTITIONS WITH VIRTUAL NODES

Below we detail the procedure to partition the full graph \mathcal{G} , described by the set of vertices \mathcal{V} and edges \mathcal{E} , into a set of slices $\{\mathcal{G}_1 \ldots \mathcal{G}_N\}$ which are augmented by virtual nodes and edges.

Algorithm 1: Augmented partitioning approach

² for $i \leftarrow 1$ to N do $\mathcal{V}_i \leftarrow [];$ $n_{i}^{n} = 0;$ for $v \in \mathcal{V}$ do if $v.x \in |x_i, x_{i+1}|$ then \mathcal{V}_i .append(v); $n_i^n += 1;$ end end $\mathcal{E}_i \leftarrow [];$ $n_{i}^{e} = 0;$ for $v_1 \in \mathcal{V}_i$ do for $v_2 \in \mathcal{V}_i$ do if $v_1 \rightarrow v_2 \in \mathcal{E}$ then \mathcal{E}_i .append $(v_1 \rightarrow v_2)$; $n_i^e += 1;$ end end end for $v_1 \in \mathcal{V}_i$ do for $v_2 \in \mathcal{V} \setminus \mathcal{V}_i$ do if $v_2 \to v_1 \in \mathcal{E}$ then \mathcal{V}_i .append(v_2); \mathcal{E}_i .append $(v_2 \rightarrow v_1)$; end end

29 | 30 end

end

 $\mathcal{G}_i(\mathcal{V}_i, \mathcal{E}_i);$

The number of labeled nodes (n_i^n) and the number of labeled edges (n_i^e) are collected and passed to the training functions, which then mask the remainder of the outputs (the virtual nodes and edge outputs) while computing the loss.

D AUGMENTED PARTITION APPROACH WITH TWO MP LAYERS

t_{slice} [Å]	N_t	# Epochs	$\epsilon_{\mathbf{node}}$	σ_{node}	ϵ_{edge}	$\sigma_{\rm edge}$
~ 2	27	16,437	8.24	21.35	0.24	0.39
~ 3	28	13,461	8.60	21.01	0.24	0.43
$\sim \!\! 4$	14	14,089	6.34	15.06	0.24	0.38
~ 6	9	11,670	4.46	11.50	0.24	0.39

Table 7: Ablation studies exploring the impact of slice length t_{slice} on the prediction accuracy with 2 MP layers. Values are reported in $\epsilon_n [mE_h]$ and $\sigma_n [\mu E_h]$.

Using 2 MP layers rather than one heavily degrades ϵ_{node} (**Table 3**). This occurs because each slice is augmented with only one layer of virtual connections, so the network does not have a correct 1026 representation of the 2-hop neighborhood. Propagation of this incorrect information into the graph 1027 thus results in a degraded ϵ_{node} , which is more sensitive to the extended environment that ϵ_{edge} . 1028 The error in this case minimally depends on the outgoing edge (albeit uniformly high). We note that 1029 the error with 2 MP layers shows a dependence on t_{slice} (**Table 9**) since larger slices require fewer 1030 virtual connections outside the partition.

- 1031
- 1032 1033

E HYPERPARAMETERS

1034 1035

We use the hyperparameters shown in Table 8 to train the H_2O benchmark and custom HfO_2 dataset. The ReduceLRonPlaeau scheduler decreases the learning rate by the decay factor when it does not detect a further decrease in validation loss within the decay patience $t_{patience}$. The threshold refers to the sensitivity of the scheduler to changes in validation loss. Once the minimum learning rate is reached, the training stops. The meaning of the hyperparameters are explained in Appendix A. Note that weight decay is not implemented in our study for all cases.

1042

Hyper-parameters	HfO ₂ /PtGe dataset	MD17 dataset
Optimizer	Adam	Adam
Precision	single (f32)	double (f64)
Scheduler	ReduceLROnPlateau	ReduceLROnPlateau
Initial learning rate	1×10^{-4}	1×10^{-4}
Minimum learning rate	1×10^{-5}	1×10^{-10}
Decay patience t _{patience}	500	50
Decay factor	0.5	0.5
Threshold	1×10^{-3}	1×10^{-5}
Interaction distance (Å)	8.0	_
Maximum degree L_{max}	4	4
Maximum order M_{max}	4	4
Number of Message Passing Layers	1	2
Number of spherical channels d_{sphere}	16	64
$f_{ii}^{(L)}$ dimension d_{attn_hidden}	16	64
Number of attention heads N_h	2	2
$f_{ii}^{(0)}$ dimension d_{attn_alpha}	16	32
Value dimension d_{attn_value}	16	32
Hidden dimension in feed forward networks d_{ffn}	64	64

Table 8: Hyper-parameters used for HfO2, PtGe and MD17 data.

1065 1066

1064

1067 1068

1069 E.1 Hyperparameter study

1070

We conducted hyperparameter tuning by training using the full training structure (structure 2), divided into 18 slices of ~ 3 Å thickness. The validation slice used for the scheduler during training was taken from the validation structure, and is centered at the location 26.5 Å. The rest of the hyperparameters follow the values in Table 8. Note that since we are tuning the hyperparameters, we evaluate the trained models with different hyperparameters on the full validation structure (structure 1). The test structure remains unseen throughout this process.

From the table, it is clear that the hyperparameters d_{sphere} , $d_{attn.value}$, $d_{attn.hidden}$ and $d_{attn.alpha}$ can all be reduced to 16 with little tradeoff in accuracy. Cutting down on the number of parameters allows us to drastically minimise the memory consumption, allowing large graphs to fit into GPU memory during training.

$d_{sphere,} \ d_{attn_hidden}$	$d_{attn_value,} \ d_{attn_alpha}$	Parameters	Epochs	$\epsilon_{\mathbf{n}}[mE_h]$	$\sigma_{\mathbf{n}}$	$\epsilon_{\mathbf{e}}[mE_h]$	$\sigma_{\mathbf{e}}$
4	32	52,572	31,940	2.82	5.72	0.36	0.84
8	32	125,324	29,002	2.54	5.73	0.18	0.35
16	16	273,644	19,106	2.39	5.31	0.16	0.33
16	32	335,436	16,857	2.46	5.84	0.19	0.35
32	32	1,014,092	14,833	2.51	5.83	0.19	0.37
64	32	3,405,132	11,000	2.47	5.83	0.19	0.37

Table 9: Table of Hyperparameters, tested on the validation structure (structure 1)

F AMORPHOUS HAFNIUM DIOXIDE (A-HFO₂) TRAINING DATA

1095 F.1 DATASET GENERATION

Atomic structures corresponding to materials in the amorphous phase can be produced through melt quench Urquiza et al. (2021), seed-and-coordinate Youn et al. (2014), or 'decorate and relax' Tafen
 & Drabold (2003) approaches. To accurately reproduce long-range structural disorder, the structures used must be large enough to avoid the creation of wavefunctions which repeat over periodic
 boundaries.

1102 The first step to computing the electronic properties of $a-HfO_2$ is to generate the atomic structure 1103 of the material in its amorphous phase. To accurately capture the structural motifs underlying this 1104 phase and a realistic range of atomic coordination (eg, neighboring Oxygens for each Hafnium, and vice versa), we start from the crystalline m-HfO₂ phase and perform melt-quench processes using 1105 Molecular Dynamics (MD), following a similar procedure as the ones described in Refs. Kaniselvan 1106 et al. (2023); Urquiza et al. (2021). We generate 3 independent structures of a-HfO₂ using the 1107 QuantumATK toolkit Søren Smidstrup et al. (2020). As the first step, we run an NVT simulation 1108 with the Langevin thermostat at 3000 K for 50 ps with a step size of 1 fs. We use the MTP-HfO₂-1109 2022 potential, provided by the software. Next, we run an NPT simulation for 300 ps (and the same 1110 1 fs step size), with an initial reservoir temperature of 3000K and a final temperature of 300K, for a 1111 cooling rate of 9K/ps. Finally, we anneal structure evolve at 300K for 50 ps, using the same NVT 1112 Langevin thermostat as for the melting. 1113

We then perform a structural relaxation with CP2K code Kühne et al. (2020) to correct for any dis-1114 crepancies between the relaxed bond lengths attained with the force field used for MD, and those 1115 obtained with DFT. Due to the computational cost of using a more complete DZVP basis set Vande-1116 Vondele & Hutter (2007), we use a simpler SZV basis VandeVondele & Hutter (2007) which uses 4 1117 basis functions per Oxygen atom and 10 basis functions per Hafnium atom. The plane-wave cutoff 1118 is set to 500 Ry, while a cutoff of 60 Ry is used for mapping the Gaussian-type orbitals onto the grid. 1119 We use the PBE functional for the exchange-correlation energy Perdew et al. (1997). To accurately 1120 capture the band gap of a-HfO₂, we applied the Hubbard correction Anisimov et al. (1991) of U = 71121 eV to the 3d orbital of Ti and the Hubbard correction of U = 10 eV to the 2p orbital of O.

1122

1080

1082 1083 1084

1086 1087 1088

1089 1090 1091

1093 1094

1123

1124 F.2 Atomic bonding environments in the amorphous phase 1125

1126 In Fig. 9 we plot the O-coordination of each Hf atom and the radial distribution function q(r) (where 1127 r is inter-atomic distance) for each of the three structures. The distribution in the coordination and 1128 dispersion of the peaks in q(r) indicates the amorphous nature of the three structures. Variations 1129 between them appear as perturbations in these two quantities. To gain more insights on how different 1130 are the structures, we additionally plot the spatially-resolved O-coordination of Hf atoms along 1131 the longest, x coordinate for the three structures, as well as the distributions of outliers (Hf atoms with very low and very high O-coordination) in three-dimensional space. It is evident that these 1132 outliers are situated at different locations in different structures, demonstrating a significant degree 1133 of dissimilarity among the structures.



Figure 9: a) O-coordination of Hf atoms (number of O atoms bonding a Hf atom) for each of the generated structures, showing a distribution around a coordination number of 6, and variation between the structures. b) The radial distribution function $(g(r) = \frac{dn(r)}{dr} \frac{V_{domain}}{4\pi r^2 N_{atoms}})$, where n(r) is the number of atoms with distance r between them for the structures 1-3. (c) Spatial distribution of coordination outliers (Hf atoms with Ocoordination equal to 8 or 4) for the three structures, which are an indicator of the uniqueness of the three structures.

1161 F.3 ENERGY EIGENVALUES 1162

1163 Due to the large cell sizes of the a-HfO₂ structures, all necessary energetic information is contained within the Γ point (where the wavevectors $k_x = k_y = k_z = 0$). The energies at this location can be 1164 computed by directly diagonalizing \mathbf{H} . Amorphous a-HfO₂ structures corresponding to a realistic 1165 distribution of bond lengths should then produce an energy bandgap. We show the distribution of 1166 energy eigenvalues for the three structures in **Table 1** in **Fig. 10**, at different values of r_{cut} . In the 1167 second row, we zoom into the range of eigenvalues around the energy bandgap, which is defined by 1168 the transition between occupied and unoccupied electronic states (the Fermi level $E_F = \sim 0.3 E_h$ in 1169 all cases). Values of $r_{cut} \ge 8\text{\AA}$ create no noticeable difference on the eigenvalue spectra. Note that 1170 the value of $r_{cut} = \infty$ corresponds to the case where no nonzero values were filtered from \mathbf{H}^{GT} . 1171

1172 1173

1174

G COMPUTE ENVIRONMENT AND RUNTIME COMPARISONS

The training is performed with PyTorch Distributed Data Parallel (Li et al. (2020)), where the graph partitions (slices) can be distributed between GPUs.

1177 1178

G.1 MEMORY CONSUMPTION OF THE FULL GRAPH

During the training of the full graph model, the peak memory consumption observed was 61.68 GiB
 on a single NVIDIA A100 GPU. Most of the consumption does not stem from the network and the structure but from the additional memory needed for the convolution operations.

1182 1183 1184

G.2 COMPUTATIONAL IMPROVEMENT WITH PARTITIONING APPROACH

In Fig. 11, we show the decrease in time per epoch and resulting speedup when using the *augmented partitioning* approach introduced in Section 3.2. Since the partitions are independent, the only communication involved in every epoch is a collective to inform each GPU/rank of the loss of each other rank. The time per epoch thus decreases uniformly with the number of slices (N_t) used.



Figure 10: Eigenvalues of the ground-truth Hamiltonian matrix, showing (left) the full eigenvalue spectrum and (right three) zoomed in around the bandgap, for the three structures, in the order of their appearance in **Table 1**.

Despite the independence of each batch and the minimal communication per epoch, the scaling is not perfectly linear. The deviation from an ideal speedup can be attributed to two factors:

• Load imbalance: The partitioning approach was designed to leverage the periodicity in the y- and z- direction within a straightforward implementation. However, it is not ideal in terms of the number of cuts/number of virtual nodes/edges required, resulting in a slightly different amount of work per rank which leads to an observable load imbalance at higher N_t . This effect can be seen in the allocated memory per partition (**Fig. 11(c)**). We note that the *augmented partitioning* method can be used with any standard graph-partitioning algorithm.

- Computational overhead of the virtual nodes and edges: Individual nodes and edges of the graph can be repeated in labeled and virtual node lists. Treating the replicas introduces additional computational cost while training the network, which increases with N_t . This overhead is maximum with the use of very small slices (large N_t), thus introducing a trade-off between parallelism and time per epoch.

1228 G.3 H_2O vs HFO_2 runtimes

In Section 4.5, we make a comparison between the computational cost of computing the Hamil-tonian for an H_2O molecule and the HfO_2 structure. However, the Hamiltonians for the H_2O molecules in the MD17 dataset were computed by Schütt et al. (2019) using different compute infrastructure. To approximate the cost of generating such a dataset under the same computational conditions, we set up CP2K simulations with a Double- ζ Valence Polarized (DZVP) basis, which in-cludes a similar set of valence and polarization functions as the def2_SVP basis used with the ORCA code (Neese (2011)) to generate the H_2O Hamiltonians for the MD17 dataset. A minor difference is that the def2_SVP basis includes an extra s-type orbital on Oxygen (14 total per Oxygen atom). Under these conditions, the computation time per H_2O molecule was 7s, when run on 12 nodes with 12-core Intel Xeon E5-2680 CPUs and NVIDIA P100 GPU (the Piz Daint supercomputer), resulting in a total of 0.04 node hours. The HfO_2 structures requires 3.65 node hours in the same compute environment (but distributed to 27 nodes). The difference, omitting scaling behavior, is roughly $\sim 100 \times$.



1254 Figure 11: (a) Time per epoch and (b) speedup resulting from the use of increasing numbers of slices N_t . 1255 Median values are shown, while the error bands are one standard deviation. Experiments were run on NVIDIA 1256 A100 GPUs with # ranks set to N_t . Measurements are only shown up to 8 slices/8 GPUs due to limitations in available compute resources at the time of submission. The fill-between indicates the range in runtime over the 1257 first 30 minutes of training. The dashed black line corresponds to the ideal speedup, in which case the use of 1258 N_t slices would enable an $N_t \times$ speedup in the runtime per epoch. (c) Measured peak memory consumption 1259 as a function of the number of partitions, where each bar corresponds to a different GPU. Variation in memory 1260 consumption between GPUs at each individual value of N_t translates to load imbalance, which correlates with 1261 the deviation from ideal scaling shown in (b).

1290 1291

1293 1294 1295

H COMPARISONS BETWEEN DFT AND PREDICTED HAMILTONIANS

24

1265 1266 In **Figure 12**, we plot the MAE error as a function of different 1267 interactions between the orbital basis of the a-HfO₂ test struc-1268 ture. The data is plotted in log scale to magnify the asymmetry 1269 resulting from the separate training of the two-way edges between labeled nodes in the graph.

1270 In Fig. 13, we plot the comparison between H^{pred} and H^{GT} 1271 (as shown in the main text) for three separate cases: (1) us-1272 ing the upper triangle, (2) the lower triangle, and (3) the sym-1273 metrized \hat{H}^{pred} . We also zoom in around the bandgap in the 1274 second row. In all three cases, H^{pred} is unchanged, indicating 1275 that the small asymmetry in the matrices caused by the exis-1276 tence of separate forward/backward edges between atoms has 1277 a minimal effect. In the third row, we show the error in log 1278 scale as a function of the eigenvalue index, sorted in the same order as shown in the first row. The error is largest around the 1279 band edges, where orbital coupling is most significant. 1280



Figure 12: Average ϵ_{MAE} between $(H_{ij})^{GT}$ and $(H_{ij})^{pred}$ for specific inter-atomic orbital coupling. The data is shown in log scale to magnify asymmetry in the error.



Figure 13: Comparison between H^{GT} and H^{pred} , using the upper/lower triangle or symmeterized version of H^{pred} . The center row shows the same plots zoomed in around the bandgap. The last row shows the difference in the eigenvalue spectra. The structures appear in the order of **Table 1**.

50	Structures		Epochs		σ [μ E]	$\int [m F]$	$\sigma \left[u F \right]$	$c_{m} [m F]$	c. [moV]
	Training	Testing		$\epsilon_{\mathbf{n}} \left[m E_{h} \right]$	$\sigma_{\mathbf{n}} \left[\mu \mathbf{L}_{h} \right]$	$\epsilon_{\mathbf{e}} \left[\mathcal{I} \mathcal{I} \mathcal{L}_h \right]$	$\sigma_{\mathbf{e}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{tot}} [\mathcal{M} \mathcal{L}_h]$	etot [mev]
	0	3	15675	2.29	5.09	0.20	0.38	0.22	5.87
	0, 2	3	48356	2.22	4.99	0.12	0.25	0.13	3.55

Table 10: Comparison of models trained on one HfO_2 structure (structure 2) vs two HfO_2 structures (0 and 2). Both are tested on a fully unseen structure (3), showing improved accuracy from 5.87 to 3.55 meV.

Oxygen v Training set	acancies Testing set	$\epsilon_{\mathbf{n}}\left[mE_{h}\right]$	$\sigma_{\mathbf{n}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{e}}\left[mE_{h}\right]$	$\sigma_{\mathbf{e}}\left[\mu E_{h}\right]$
5%	5%	2.44	5.06	0.16	0.31
5%	10%	2.58	5.23	0.18	0.33
5%	15%	2.50	4.96	0.17	0.33
10%	5%	2.48	5.12	0.18	0.33
10%	10%	2.50	4.96	0.17	0.33
10%	15%	2.60	4.89	0.18	0.33
15%	5%	2.94	6.45	0.16	0.34
15%	10%	2.52	5.01	0.16	0.34
15%	15%	2.52	4.69	0.16	0.31

Table 11: HfO₂ models trained and tested with different stoichiometry using augmented partitioning. 18 slices of each structure, each 3 Å thick, was used for training. The training method is identical to the one used to obtain Table 4. Models trained on vacancies 5%, 10 % and 15% vacancies are tested on test structures with vacancies ranging from 5% to 15%.

1372 1373

1375

1377

1363 1364 1365

1367 1368

1374 I ADDITIONAL TESTS

1376 I.1 USING TWO HFO₂ STRUCTURES FOR TRAINING

While training on a single structure is sufficient to achieve the results in **Fig. 6**, increasing the quantity of training data further leads to improved prediction accuracy. To demonstrate this, we train using *augmented partitioning* on two structures, including structure 2 from **Table 1** as well as a newly generated structure "0" of similar size (3000 atoms). Training was performed with 18 slices of 3 Å thickness taken from both structures, validated on structure 1, and tested on structure 3. The results are compared against the previous model trained on one structure in Table 10, showing that the prediction accuracy can be improved from 5.87 *meV* to 3.55 *meV* by using 2× the number of slices for training.

1385

1386 I.2 SUB-STOICHIOMETRIC HAFNIUM OXIDE

Amorphous HfO_2 often exists in a sub-stoichiometric form (HfO_x) , which can be interpreted as the presence of oxygen vacancies. In this section we evaluate whether our model can extend to train and predict such defective structures.

1391 To do this, we create a dataset for sub-stoichiometric HfO_x structures by introducing randomly 1392 distributed oxygen vacancies into the original, pristine HfO_2 structures. The sub-stoichiometric 1393 structures are generated for x = 1.9, 1.8, and 1.7 (corresponding to vacancy concentrations of 5%, 1394 10%, and 15 %, respectively). Vacancies are treated as ghost atoms (atoms with no orbitals, but with 1395 a basis set defined at their locations), to mitigate the basis set superposition error Senent & Wilson (2001), a known problem related to localized basis sets. More precisely, by treating vacancies as 1396 ghost atoms, one prevents the excessive borrowing of the basis sets from neighboring atoms by the 1397 vacancy, which improves the accuracy of the predicted electronic properties. These ghost atoms are 1398 assigned an atomic number of 0. The training and testing approach is similar to the one used to 1399 obtain Table 4, except that now oxygen vacancies are considered. For all experiments, 18 slices (3 1400 Athick) were used. 1401

1402 The results are summarized in Table 11. The ϵ_n and ϵ_e values across different experiments lie within 1403 a small range (2.50-2.90 mE_h and 0.16-10.18 mE_h respectively), showing that the network generalizes well to structures of different vacancies, regardless of which vacancy configuration it was

1404 1405	Training method method	Oxygen vacancies Testing set	$\epsilon_{\mathbf{n}}\left[mE_{h}\right]$	$\sigma_{\mathbf{n}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{e}}\left[mE_{h}\right]$	$\sigma_{\mathbf{e}}\left[\mu E_{h}\right]$
1406	partitioned	5%	2.94	6.45	0.16	0.34
1407	partitioned	10%	2.52	5.01	0.16	0.34
1408	partitioned	15%	2.52	4.69	0.16	0.31
1409	full	5%	2.96	6.07	0.19	0.38
1410	full	10%	2.67	5.18	0.18	0.36
1411	full	15%	2.64	4.83	0.17	0.35

Table 12: Comparison between full graph training and the augmented partitioning training using the same HfO_2 structure with 15% vacancies. Models are tested on structures with vacancies ranging from 5% to 15%.

Material	Cutoff [Å]	$\epsilon_{\mathbf{n}}\left[mE_{h}\right]$	$\sigma_{\mathbf{n}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{e}}\left[mE_{h}\right]$	$\sigma_{\mathbf{e}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{tot}}\left[mE_{h}\right]$	$\epsilon_{\mathbf{tot}} \left[meV \right]$
Crystalline HfO ₂	8	0.01	0.02	0.04	0.07	0.04	1.17
Amorphous HfO ₂	8	2.29	5.09	0.20	0.36	0.22	5.87
Amorphous PtGe	16	0.87	1.43	0.05	0.10	0.05	1.43

1420 Table 13: Summary of models trained on crystalline HfO₂, amorphous HfO₂ and amorphous PtGe materials, respectively. The HfO_2 model was trained on different slices of the same crystalline structures. On the other 1422 hand, the PtGe model was trained on a single 5 Å slice of structure 1 and tested on on a fully unseen structure 1423 2.

trained on. To demonstrate that the *augmented partitioning* approach similarly does not affect ac-1426 curacy for sub-stoichiometric HfO_x, we also perform full graph training using structure 2 with 15%1427 vacancies, and compare with the augmented partitioning approach in Table 12. The minimal dif-1428 ference in ϵ_n and ϵ_e values between full and partitioned approaches indicates that both approaches 1429 generalize equally well to different stoichiometry. These values are also close to that of stoichio-1430 metric HfO_2 in Table 4, demonstrating that the augmented partitioning approach can also be applied 1431 even in the case of more realistic sub-stoichiometric structures.

1433 I.3 CRYSTALLINE HAFNIUM OXIDE 1434

Crystalline materials contain highly regular atomic environments, and are thus a natural extension of 1435 this approach. Although such a large unit cell is not required for a crystal, we nevertheless generate 1436 a single crystalline HfO_2 structure (in its monoclinic phase) containing 3000 atoms for comparison. 1437 The model was trained on a single slice taken from x = 0 Å, validated on a slice from x = 15 Å, 1438 and tested on an unseen slice from x = 20 Å. Results in Table 13 show that a high accuracy close to 1439 sub-meV can be achieved. 1440

I.4 SECOND MATERIAL EXAMPLE: PTGE 1442

1443 Finally, we test if the model and training approach can be used for different material systems. We 1444 consider the example of two amorphous PtGe structures, labelled 1 and 2, each containing 2688 1445 atoms. These structures were similarly generated in CP2K. A larger cutoff radius of 16 Å was 1446 chosen due to the larger spacing between atoms. The model was trained on a single 5 Å slice taken 1447 from x = 10 Å and validated on another slice at x = 20 Å, both from structure 1. The trained model 1448 was then tested on a full unseen structure (structure 2), with results shown in Table 13. The final obtained error for all 2688 atoms and 2,148,055 edges is 1.43 meV, much lower than that of HfO₂, 1449 despite the lower amount of training data. This demonstrates the generalizability and robustness of 1450

1451 1452 1453

1412

1413

1414

1421

1424 1425

1432

1441

Cutoff $[Å]$	$\epsilon_{\mathbf{n}}\left[mE_{h}\right]$	$\sigma_{\mathbf{n}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{e}}\left[mE_{h}\right]$	$\sigma_{\mathbf{e}}\left[\mu E_{h}\right]$	$\epsilon_{\mathbf{tot}} \left[m E_h \right]$	$\epsilon_{tot} [meV]$
6	0.87	1.43	0.15	0.19	0.17	4.60
8	0.87	1.42	0.09	0.16	0.10	2.73
16	0.87	1.43	0.05	0.10	0.05	1.43

Table 14: Prediction accuracy of model on amorphous PtGe material with different r_{cut}

the augmented partitioning approach when applied to different material dataset with a much larger cutoff radius.

We perform a similar study on the cutoff radius and connectivity of the PtGe, through both eigenvalue analysis and the convergence study of cutoff-radii, with results shown in Table 14 and Figure 14. Increasing the cutoff radius once again increases the overall prediction accuracy of the trained model, with the improvement especially noticeable at the edges.



Figure 14: Atomic coordinates of the two amorphous PtGe structures.

I THEORETICAL JUSTIFICATION FOR AUGMENTED PARTITIONING APPROACH

The augmented partitioning approach relies on the correction of the partitioned sub-graphs through
the introduction of virtual nodes and edges, such that partitioned training fully resembles full graph
training. Here, we provide the theoretical foundations upon which our approach is built.

We start by recalling the operation of a message passing neural network trained on a full graph. During training, the nodes are updated through the aggregation of messages from connected neighbors. More specifically, in the first layer of equivariant message passing networks, the n'_i inputs to the aggregation function are simply the atomic numbers Z_i and Z_j (embedded in a tensor) of atom *i* and neighbor *j*, their scalar distances $|r_{ji}|$ (expanded in the Gaussian basis), and the normalized vector \hat{r}_{ji} indicating the orientation of the edges (embedded within the rotation and counter-rotation operations)

$$n'_{i} = \sum_{j \in \mathcal{N}(i)} \Phi_{j}(Z_{i}, Z_{j}, |r_{ji}|, \hat{r}_{ji}).$$
(1)

In Eq. (1), $\mathcal{N}(i)$ is the neighbor list of atom i and Φ_i is a learnable function that encompasses all the operations of our equivariant network (convolution, gate activation, attention weights). The sum of these functions over all neighbors represents the overall aggregation function that we aim to learn using our equivariant network. It maps the inputs to the output node embedding n_i '. Similarly, for edges, the updated node embeddings e'_{ji} fed into a learnable function Φ_{ji} has the following form:

1500

1490 1491

1492 1493

1473

1474 1475 1476

1479

$$'_{ji} = \Phi_{ji}(n'_i, n'_j, |r_{ji}|, \hat{r}_{ji}).$$
 (2)

1503 It maps the inputs consisting of the updated node embeddings to the edge embeddings. In our case, 1504 due to the large unfeasible size of the graph, we have to partition it into slices for training. In 1505 partitioned subgraphs, however, there are also connected neighbor nodes that lie outside of the par-1506 tition, meaning that some of the j terms in Eq. (1) are missing. Ignoring the contribution from these 1507 nodes leads to incomplete/wrong aggregation. As a result, the wrong aggregation function would 1508 be learned when fitting the final node output to the target Hamiltonian data (during minimization of 1509 MSE loss). See the ablation study in Section 4 and Table 3) for more details.

e

This is why we introduced virtual nodes and edges. They account for the presence of connected
 neighbors outside of the partition when computing the updated node embeddings for atoms situated within the partition:



Figure 15: Scatter plots showing the decay of MAE with increasing distance from different perturbations, including (a) 0.1 Å translation of an O atom, (b) replacement of O atom with a vacancy, and (c) replacement of O atom with Hf atom

 $n'_{i\in\mathcal{P}} = \sum_{j\in\mathcal{N}(i)\cap j\in\mathcal{P}} \Phi_j(Z_i, Z_j, |r_{ji}|, \hat{r}_{ji}) + \sum_{j\in\mathcal{N}(i)\cap j\notin\mathcal{P}} \Phi_j(Z_i, Z_j, |r_{ji}|, \hat{r}_{ji}),$

1525

1527

1520

1530 1531

$$e'_{(i,j)\in\mathcal{P}} = \Phi_{ji}(n'_{i\in\mathcal{P}}, n'_{j\in\mathcal{P}}, |r_{ji}|, \hat{r}_{ji}), \tag{4}$$

(3)

where \mathcal{P} is defined as the set of atoms belonging to the partition. The virtual nodes and edges are 1532 cast into the second summation on the right-hand-side of Eq. (3). They contain all necessary inputs 1533 $(Z_i, Z_i, |r_{ii}|, \hat{r}_{ii})$ needed to compute the aggregation function in the first layer of the network. 1534 Therefore, for the case of nodes within the partition $(n_{i \in \mathcal{P}})$, Eq. (3) is now equivalent to Eq. (1), 1535 as they contain the same terms and inputs. By extension, since the edge embeddings within the 1536 partition are only updated based on nodes within the partition (Eq. (4)), they are also correct, with 1537 Eq. (4) being equivalent to Eq. (2) for all $e_{(i,j)\in\mathcal{P}}$. Overall, in the case of a single MP layer, the local 1538 environment for nodes and edges within the partition is thus identical to that of the full graph. As 1539 a consequence, the correct aggregation function, along with accurate predictions, are obtained from 1540 training.

Note that throughout this process, the output embeddings of the virtual nodes and edges are not used at all and remain completely masked during training - only their inputs $(Z_i, Z_j, |r_{ji}|, \hat{r}_{ji})$ are used to inform the network. Their own local environment and outputs have, therefore, no influence on the aggregation function learned, and do not need to be corrected.

1546

1547

1548 1549

K ANALYSIS OF PERTURBATION EFFECTS AT VARYING DISTANCES

To demonstrate the effects of long and short range perturbations in amorphous structures, we introduce a single perturbation at one chosen location in the structure and measure the mean absolute error of the onsite Hamiltonian blocks when compared to that of the unperturbed structure. The types of perturbations introduced include single oxygen atom translation, oxygen vacancy, and substitution of an oxygen by a hafnium atom, which are plotted against distance from perturbation in Fig. 15 (a), (b) and (c) respectively. In all cases, the effect of the perturbation rapidly decays with increasing distance.

For the case of a 0.1 Å translation perturbation, the average onsite MAE at a distance of 8 Å away is given by 0.15 mE_h . Considering the average value of an onsite Hamiltonian block (63 mE_h), the perturbation only affects the matrix elements by 0.24% overall. Similarly, for vacancy and substitution perturbations, the matrix elements of atoms located 8 Å away only changed by 0.18% and 0.12% respectively. This implies that for our chosen cutoff of 8 Å, perturbations occurring outside of the radius surrounding the atom have a negligible effect on its Hamiltonian matrix elements. This also means that the electronic structure of that atom can be learned using information from the local atomic environment.

This is also demonstrated through our study of sub-stoichiometric HfO_x with randomly distributed vacancies. Despite training on independent slices, we ensure that every atom within that slice is

1566	surrounded by a complete local environment with a radius of 8 Å through the use of virtual nodes
1567	and edges. Any vacancies outside of that radius have a negligible effect on the atom, and are seen
1568	and learned by other partitions. When multiple slices are trained together, the entire distribution
1569	of perturbations are captured, allowing the model to generalize well to unseen structures with a
1570	completely different distribution of vacancies and local atomic environments.
1571	
1572	
1573	
1574	
1575	
1576	
1577	
1570	
1579	
1500	
1501	
1582	
1584	
1585	
1586	
1587	
1588	
1589	
1590	
1591	
1592	
1593	
1594	
1595	
1596	
1597	
1598	
1599	
1600	
1601	
1602	
1603	
1604	
1605	
1606	
1607	
1608	
1609	
1610	
1611	
1612	
1613	
1014	
1616	
1617	
1612	
1619	
1010	